

(11)特許出願公開番号
特開2002-324000
(P2002-324000A)

(43)公開日 平成14年11月8日(2002.11.8)

(51) Int.CL ⁷		識別記号	F I	データコード [*] (参考)	
G 0 6 F	12/00	5 1 4	G 0 6 F 12/00	5 1 4 E	5 B 0 6 5
		5 3 1		5 3 1 D	5 B 0 8 2
		5 4 5		5 4 5 A	
	3/06	3 0 4	3/06	3 0 4 F	

審査請求 未請求 請求項の数17 OL 外国語出願 (全 54 頁)

(21) 出願番号	特願2002-3193(P2002-3193)	(71) 出願人	000005108 株式会社日立製作所 東京都千代田区神田駿河台四丁目 6 番地
(22) 出願日	平成14年 1 月10日 (2002. 1. 10)	(72) 発明者	渡辺 直企 アメリカ合衆国 カリフォルニア州 プリ スバーン シェアラポイントパークウェイ 2000 日立アメリカ リサーチアンドディ ベロップメント部内
(31) 優先権主張番号	0 9 / 7 6 0 3 4 4	(72) 発明者	山本 彰 神奈川県川崎市麻生区王禅寺1099番地 株 式会社日立製作所システム開発研究所内
(32) 優先日	平成13年 1 月12日 (2001. 1. 12)	(74) 代理人	100091096 弁理士 平木 祐輔
(33) 優先権主張国	米国 (U.S.)		

最終頁に続く

[最終頁に続く](#)

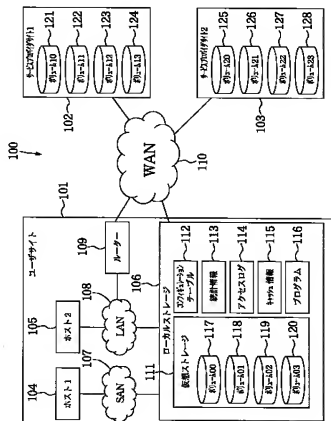
(54) 【発明の名称】 仮想ボリューム・ストレージ

(57) 【要約】

【課題】 データ・ストレージ・サービスを提供するシステムを提供する。

【解決手段】 データ・ストレージ・サービスを提供するように構成されているサービス・プロバイダ・サイト 102、103 及びサービス・プロバイダ・サイトにて W A N 110 により結合されたユーザ・サイト 101 を備える。ユーザ・サイトは仮想ストレージ 111 を持つローカル・ストレージ 106 を備え、仮想ストレージは同期ボリューム及び非同期ボリュームを備え、ローカル・ストレージは同期ボリューム内に書き込まれたデータをサービス・プロバイダ・サイトに即座に送信し、非同期ボリュームに書き込まれたデータを所定のスケジュールでサービス・プロバイダ・サイトに送信し、データがローカル・ストレージに格納されていない場合にサービス・プロバイダ・サイトからデータを読み込むように構成されている。

图 1



【特許請求の範囲】

【請求項1】 データ・ストレージ・サービスを提供するシステムであって、

データ・ストレージ・サービスを提供するように構成されているサービス・プロバイダと、

ワイド・エリア・ネットワーク（WAN）によりサービス・プロバイダ・サイトに結合されたユーザ・サイトを備え、前記ユーザ・サイトは仮想ストレージを持つローカル・ストレージを備え、前記仮想ストレージは同期ボリューム及び非同期ボリュームを備え、前記ローカル・ストレージは前記同期ボリューム内に書き込まれたデータを前記サービス・プロバイダ・サイトに即座に送信し、前記非同期ボリュームに書き込まれたデータを所定のスケジュールで前記サービス・プロバイダ・サイトに送信し、データが前記ローカル・ストレージに格納されていない場合に前記サービス・プロバイダ・サイトからデータを読み込むように構成されているシステム。

【請求項2】 前記ローカル・ストレージがホストによる当該ローカル・ストレージへのアクセス活動を記録するように構成されていることを特徴とする請求項1記載のシステム。

【請求項3】 前記ローカル・ストレージが前記仮想ストレージ内のボリュームへのアクセス活動を記録することとを特徴とする請求項1記載のシステム。

【請求項4】 前記ローカル・ストレージが前記仮想ストレージ内のボリュームのサブエリアでのアクセス活動を記録することを特徴とする請求項1記載のシステム。

【請求項5】 分析のため前記サービス・プロバイダ・サイトでアクセス活動を行えることを特徴とする請求項2記載のシステム。

【請求項6】 データ・ストレージ・サービスを提供する方法であって、

同期ボリュームと非同期ボリュームとを有する仮想ストレージを備えるローカル・ストレージを提供するステップと、

選択的に実行するステップであって、

前記仮想ストレージの同期ボリューム内に書き込まれたデータをサービス・プロバイダ・サイトに即座に送信するステップと、

前記仮想ストレージの非同期ボリューム内に書き込まれたデータを前記サービス・プロバイダ・サイトに所定のスケジュールで送信するステップと、
データが前記ローカル・ストレージ内に格納されていない場合に前記サービス・プロバイダ・サイトからデータを読み込むステップのうち少なくとも1つを選択的に実行するステップとを含む方法。

【請求項7】 前記ローカル・ストレージへのアクセス活動を記録するステップを更に含むことを特徴とする請求項6記載の方法。

【請求項8】 前記サービス・プロバイダ・サイトへの

アクセス活動の記録を提供するステップを更に含むことを特徴とする請求項7記載の方法。

【請求項9】 仮想ストレージの同期ボリューム内に書き込まれたデータをサービス・プロバイダ・サイトに即座に送信するステップ、

前記仮想ストレージの非同期ボリューム内に書き込まれたデータを前記サービス・プロバイダ・サイトに所定のスケジュールで送信するステップ、及びデータが前記仮想ストレージを含むローカル・ストレージ内に格納されていない場合に前記サービス・プロバイダ・サイトからデータを読み込むステップのうち少なくとも1つを選択的に実行するステップを含むデータ・ストレージ・サービスを提供する方法をコンピュータに実行させるためのプログラムを格納する電子的に読み取り可能な媒体。

【請求項10】 仮想ストレージの同期ボリューム内に書き込まれたデータをサービス・プロバイダ・サイトに即座に送信するステップ、

前記仮想ストレージの非同期ボリューム内に書き込まれたデータを前記サービス・プロバイダ・サイトに所定のスケジュールで送信するステップ、及びデータが前記仮想ストレージを含むローカル・ストレージ内に格納されていない場合に前記サービス・プロバイダ・サイトからデータを読み込むステップのうち少なくとも1つを選択的に実行するステップを含むデータ・ストレージ・サービスを提供する方法をコンピュータに実行させるための搬送波に埋め込まれたプログラム・コード。

【請求項11】 データ・ストレージ・システム内のデータを処理する装置であって、

仮想ストレージを持つローカル・ストレージを備えるユーザ・サイトであって、前記仮想ストレージは同期領域及び非同期領域を含み、前記ローカル・ストレージは前記同期領域内に書き込まれたデータをサービス・プロバイダ・サイトに即座に送信し、前記非同期領域に書き込まれたデータを所定のスケジュールで前記サービス・プロバイダ・サイトに送信し、データが前記ローカル・ストレージに格納されていない場合に前記サービス・プロバイダ・サイトからデータを読み込むように構成されているユーザ・サイトを含むことを特徴とする装置。

【請求項12】 ローカル・ストレージへのコマンド信号に応答して、

仮想ストレージ内の同期領域に書き込まれたデータをサービス・プロバイダ・サイトに即座に送信するステップ、

前記仮想ストレージ内の非同期領域に書き込まれたデータを前記サービス・プロバイダ・サイトに所定のスケジュールで送信するステップ、

データが前記ローカル・ストレージ内に格納されていない場合に前記サービス・プロバイダ・サイトからデータを読み込むステップのうち少なくとも1つを選択的に実行するデータ・ストレージ・システム内のデータを処理

する方法。

【請求項 13】 データ・ストレージ・システム内のデータを読み込み方法であって、

ホストからローカル・ストレージへ読み込みコマンドを受信するステップと、

前記読み込みコマンドで要求されたデータのボリューム・タイプを判別するステップと、

前記ボリューム・タイプが静的である場合に、前記ローカル・ストレージ内のローカル・ディスクからデータを読み込むステップと、

前記ボリューム・タイプがキャッシュの場合に、データが前記ローカル・ディスクに格納されているかどうかを調べるステップと、

前記データが前記ローカル・ディスクに格納されている場合に、前記ローカル・ディスクからデータを読み込むステップと、

前記データが前記ローカル・ディスクに格納されていない場合に、リモート・サービス・プロバイダ・サイトからデータを読み込み、前記データをローカル・ディスクに格納し、前記ローカル・ディスクから前記データを読み込むステップと、

前記ローカル・ディスクからデータを読み込んだ後に、前記読み込みコマンドに関係する統計情報を更新するステップと、

前記データを前記ホストに返すステップとを含む方法。

【請求項 14】 データ・ストレージ・システム内のデータを処理する装置であって、

仮想ストレージ内の同期領域に書き込まれたデータを即座に受信し、前記仮想ストレージ内の非同期領域に書き込まれたデータを所定のスケジュールで受信し、データがローカル・ストレージに格納されていない場合に前記仮想ストレージを含む前記ローカル・ストレージにデータを送信するように構成されているサービス・プロバイダ・サイトを備える装置。

【請求項 15】 仮想ストレージ内の同期領域に書き込まれたデータを即座に受信するステップ、

前記仮想ストレージ内の非同期領域に書き込まれたデータを所定のスケジュールで受信するステップ、及びデータがローカル・ストレージ内に格納されていない場合に前記仮想ストレージを含む前記ローカル・ストレージにデータを送信するステップのうち 1 つを選択的に実行するデータ・ストレージ・システム内のデータを処理する方法。

【請求項 16】 データ・ストレージ・システム内にデータを書き込む方法であって、

ホストからローカル・ストレージへの書き込みコマンドを受信するステップと、

前記書き込みコマンドのデータのボリューム・タイプを判別するステップと、

前記ボリューム・タイプが静的である場合に、前記ロー

カル・ストレージ内のローカル・ディスクに前記データを書き込むステップと、

前記ボリューム・タイプがキャッシュの場合に、前記データが前記ローカル・ディスクに格納されているかどうかを調べるステップと、

前記データが前記ローカル・ディスクに格納されている場合に、キャッシュから前記ローカル・ディスクに前記データを書き込むステップと、

前記データが前記ローカル・ディスクに格納されていない場合に、前記ローカル・ディスクと前記キャッシュ内にデータ領域を割り当てて、前記データを前記キャッシュから前記ローカル・ディスクに書き込むステップと、前記データのボリューム・タイプが同期かどうかを調べるステップと、

前記ボリューム・タイプが同期の場合に、前記データをリモート・サービス・プロバイダ・サイトと即座に同期させ、前記書き込みコマンドに関係する統計情報を更新するステップと、

前記ボリューム・タイプが同期でない場合に、前記データをリモート・サービス・プロバイダ・サイトと所定のスケジュールに基づいて同期させ、前記書き込みコマンドに関係する統計情報を更新するステップとを含む方法。

【請求項 17】 データ・ストレージ・サービスを提供するシステムであって、

仮想ストレージの同期ボリューム内に書き込まれたデータをサービス・プロバイダ・サイトに即座に送信、

前記仮想ストレージの非同期ボリューム内に書き込まれたデータを前記サービス・プロバイダ・サイトに所定のスケジュールで送信、及びデータが前記仮想ストレージを含むローカル・ストレージ内に格納されていない場合にサービス・プロバイダ・サイトからのデータ読み込みのうち少なくとも 1 つを選択的に実行する手段を備えるシステム。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】本発明はデータ・ストレージ・システムに関し、より詳細には、ユーザ・サイトに仮想ボリューム・データ・ストレージを備えるシステムと方法に関する。

【0002】

【従来の技術】大規模なデータ・ストレージ・システムを管理することは非常に困難である。通常、データ・ストレージ・サービス・プロバイダは、ストレージ・ボリューム、データ・バックアップなどを実施するストレージ関連サービスを提供する。ユーザ・サイトから、ユーザがストレージ・サービス・プロバイダのディレクトリ経由でこのストレージに接続する場合、ユーザは、長距離接続を利用してこのストレージに接続する必要がある。このプロセスは、ユーザがローカル・ストレージに

接続する場合に比べて遅延時間が長くなる。

【0003】World Wide Web (WWW) は、広い地域でストレージ・システムとして効果的な働きをする。ユーザは、ユーザ・サイトにプロキシ・サーバを配備することができる。当業者には周知のように、プロキシ・サーバはクライアント・アプリケーション (Webブラウザなど) とリモート・サーバとの間に置かれるサーバである。プロキシ・サーバは、リモート・サーバ上で利用できるアイテムをキャッシュする機能を備える。プロキシ・サーバは、リモート・サーバに送られるすべての要求を横取りし、プロキシ・サーバ側でその要求を遂行できるかどうかを調べる。プロキシ・サーバ側でその要求を遂行できなければ、プロキシ・サーバはその要求をリモート・サーバに転送し、リモート・サーバに処理させる。プロキシ・サーバは、リモート・サーバ又はサイトから送られたキャッシュ・データを読み込むことができるだけであって、リモート・サーバ又はサイトへの書き込み手順をサポートしているわけではない。さらに、プロキシ・サーバでは、データの特徴に依存するサービスを提供することもできない。そのため、プロキシ・サーバでは、資源を効果的に利用することができず、またプロキシ・サーバを使用した場合、ローカル・ストレージを使用した場合に比べて遅延時間が長くなる。

【0004】米国特許第5,155,835号では、マルチレベルの階層型動的マッピング・データ・ストレージ・サブシステムを開示している。この特許では、ワイド・エリア・ネットワーク (WAN) 間のストレージ・システムを開示しておらず、アクセス・インタフェース・タイプ (ブロック又はファイルといったもの)、目的 (ユーザ・データ又はシステム・ファイル又はバックアップ・ファイルといったもの) などのデータの特徴を無視している。したがって、この引用で開示されているサブシステムは不効率である。

【0005】インターネット・プロトコル小型コンピュータ周辺機器インタフェース (iSCSI) では、ネットワーク・インフラにインターネット・プロトコル (IP) を利用し、既存のローカル・エリア・ネットワーク及び/又はワイド・エリア・ネットワーク上で大量のブロック・ストレージ (SCSI) データを素早く転送する。iSCSI (IP SAN) では、すべての主要ネットワーク・プロトコルをサポートできるため、企業全体のネットワーク・アーキテクチャを統一でき、それにより、全体のネットワーク・コスト及び複雑度を低減できる。信頼性を保証するために、iSCSI ではIPネットワーク用に開発された周知のネットワーク管理ツール及びユーティリティを使用できる。たとえば、iSCSIプロトコルは、Webサイト (<http://www.ece.cmu.edu/ips/index.html>) のIPストレージ・セクション (<http://www.ece.cmu.edu/ips/Docs/docs.html>) で説明されている。

【0006】IETF (Internet Engineering Task Force) の作業グループは (<http://www.ietf.org>) で、ネットワーク・ファイル・システム (NFS) バージョン3 (IETFのrfc1813) 及び共通のインターネット・ファイル・システム (CIFS) プロトコル (<http://www.cifs.org>) を提供している。

【0007】NFSは、オープン・オペレーティング・システムであり、このOSを利用すればすべてのネットワーク・ユーザが異なる種類のコンピュータに格納されている共有ファイルにアクセスできる。NFSでは、伝送制御プロトコル/インターネット・プロトコル (TCP/IP) の上位で実行される仮想ファイル・システム (VFS) と呼ばれるインタフェースを通じて共有ファイルにアクセスできる。NFSを使用すると、ネットワークに接続されているコンピュータをクライアントとして動作させながらリモート・ファイルにアクセスし、サーバとして動作させながらリモート・ユーザによるローカル共有ファイルへのアクセスを可能にできる。

【0008】CIFSプロトコルは、数百万台のコンピュータに一度にリモート・ファイルにアクセスできるようにする標準を定めるものである。CIFSでは、プラットフォームもコンピュータも異なるユーザでも、新規ソフトウェアをインストールすることなくファイルを共有できる。CIFSは、TCP/IP上で実行されるが、ファイル及びプリンタ・アクセス用にMicrosoft Windows (登録商標) にあるサーバ・メッセージ・ブロック (SMB) プロトコルを使用する。したがって、CIFSでは、すべてのアプリケーション (Webブラウザを含む) がインターネット上のファイルを開き、共有することができる。

【0009】Webサイト <http://www.cnt.com> 及び <http://www.san.com> には、ストレージ又はストレージ・エリア・ネットワーク (SAN) にワイド・エリア・ネットワークを接続する方法についての説明がある。WAN接続は、非同期転送モード (ATM)、同期光ネットワーク (SONET) などが考えられる。上記の引用では、ストレージ・システムとホスト・デバイスとの間の階層型管理手法を実現する方法を開示していない。

【0010】従来のシステムと上述の欠陥を克服するシステムと方法が必要である。さらに、アクセス・サービス・プロバイダが信頼できるストレージ・システムを備えることができ、またユーザがストレージ・システムに素早くアクセスできるシステムと上述の欠陥を必要である。さらに、アクセス・サービス・プロバイダがユーザ統計データ及びユーザ・ログ・データに基づいてローカル・ストレージ・システムをチューニングできるシステムと方法も必要である。

【0011】

【発明が解決しようとする課題】本発明は、ユーザ・サイトに信頼性の高い仮想ローカル・ストレージを使いや

すい形で実現することができる。本発明はさらに、ユーザ向けの高速なアクセスを可能にする仮想ローカル・ストレージを使いやすい形で実現することができる。本発明はさらに、ローカル・ストレージで追跡されるユーザ統計データ及びユーザ・ログ・データを使用してローカル・ストレージをサービス・プロバイダ側で使いやすい形でチューンアップできるようにし、サービス・プロバイダはこのようなチューニング・サービスについてユーザに課金できる。

【0012】

【課題を解決するための手段】本発明の一実施形態では、サービス・プロバイダはワイド・エリア・ネットワーク経由でデータ・ストレージ・サービスを提供することができる。仮想ボリューム・ストレージはユーザ・サイトに配備される。ユーザ・サイトのホストは、仮想ボリューム・ストレージを使用することによりサービス・プロバイダ・サイトに配置されているサービス・プロバイダ・ストレージにアクセスできる。仮想ボリューム・ストレージを利用すると、ユーザはユーザ・サイトとサービス・プロバイダ・サイトの間で結合されているワイド・エリア・ネットワークに毎回アクセスする必要がなくなり便利である。

【0013】本発明の一実施形態では、データ・ストレージ・サービスを提供するシステムであって、データ・ストレージ・サービスを提供するように構成されているサービス・プロバイダ・サイト及びサービス・プロバイダ・サイトにワイド・エリア・ネットワーク（WAN）により結合されたユーザ・サイトを備え、ユーザ・サイトは仮想ストレージを持つローカル・ストレージを備え、仮想ストレージは同期ボリューム及び非同期ボリュームを備え、ローカル・ストレージは同期ボリューム内に書き込まれたデータをサービス・プロバイダ・サイトに即座に送信し、非同期ボリュームに書き込まれたデータを所定のスケジュールでサービス・プロバイダ・サイトに送信し、データがローカル・ストレージに格納されていない場合にサービス・プロバイダ・サイトからデータを読み込むように構成されているシステムを広範にわたって提示する。

【0014】本発明の他の実施形態では、データの特徴に基づいてローカル・ストレージ内のデータを管理する方法を提示する。ローカル・ストレージ内の仮想ボリューム・ストレージは、データが静的データかキャッシュされたデータか、同期データか非同期データを判別する。仮想ボリューム・ストレージはさらに、ストレージ・ボリュームだけでなく、ディレクトリ、ファイル、シリンダ、及び／又はブロック・アドレスでもデータを管理できる。

【0015】本発明の他の実施形態では、ユーザ・サイト内のユーザのアクセス活動を追跡する方法を提示する。ユーザ・アクセス活動を記録することで、サービス

・プロバイダはユーザ・アクセス活動のパターン分析できる。この分析に基づき、サービス・プロバイダはユーザ・サイト内の仮想ボリューム・ストレージをチューニングできる。

【0016】

【発明の実施の形態】以下では、当業者が本発明を実施し使用できるように説明し、また特定のアプリケーションとその要件の文脈において説明を行っている。実施形態のさまざまな修正は、当業者であればやすく理解できるであろうし、またここで定義している一般的な原理は本発明の精神と範囲を逸脱することなく他の実施形態及びアプリケーションにも適用できる。そこで、本発明は示されている実施形態に限定されるのではなく、ここで開示されている原理、特徴および教示に一致する最も広い範囲を許容するものとする。

【0017】本発明による一実施形態では、システム100はユーザ・サイト101、さらに第1のサービス・プロバイダ・サイト102及び第2のサービス・プロバイダ・サイト103により構成される。ユーザ・サイト101の詳細を以下で説明する。システム100内のサービス・プロバイダ・サイトの数は異なってもよい。第1のサービス・プロバイダ・サイト102は、ストレージ・ボリューム121、122、123、及び124を備え、第2のサービス・プロバイダ・サイト103はストレージ・ボリューム125、126、127、及び128を備える。サービス・プロバイダ・サイト内のストレージ・ボリュームの数は異なってもよい。サービス・プロバイダ・サイト102及び103は、たとえば、2つの異なる安全な建物内に配置できる。ユーザ・サイト101、第1のサービス・プロバイダ・サイト102、及び第2のサービス・プロバイダ・サイト103はそれぞれ、ワイド・エリア・ネットワーク（WAN）110に接続されている。サービス・プロバイダ102及び103は、WAN 110経由でデータ・ストレージ・サービスをユーザ・サイト101のユーザに提供する。

【0018】当業者には周知のとおり、WANとは、通常比較的に広い範囲にわたる地理的領域にまたがるコンピュータ・ネットワークのことである。WANはまた、ローカル・エリア・ネットワーク（LAN）も含む。WANに接続されているコンピュータは、多くの場合、電話システムなどの公衆網を通じて接続される。さらに、専用回線や衛星を通じて接続することもできる。現存最大のWANはインターネットである。

【0019】WAN 110では、各サイト（ユーザ・サイト101とサービス・プロバイダ・サイト102及び103）の間で通信ができる。各サイト101、102、及び103は、長い距離で互いに隔てられてよい。WAN 110は、通常、非同期転送モード（ATM）、同期光ネットワーク（SONET）、高密度波長

分割多重 (DWDM)、又はインターネット・プロトコル (IP) ネットワークにより構成される。

【0020】ユーザ・サイト101では、第1のホスト104はストレージ・エリア・ネットワーク (SAN) 107を介してローカル・ストレージ106に接続され、第2のホスト105はローカル・エリア・ネットワーク (LAN) 108を介してローカル・ストレージ106に接続されている。ホスト104及び105はローカル・ストレージ106にアクセスする。ホスト104及び105は、たとえば、サーバである。ユーザ・サイト101内のホストの数は異なってもよい。LAN 108は、たとえば、ルータ109を介してWAN 110に接続されている。本発明の教示から、SAN 107又はLAN 108のいずれか1つ、又はSAN 107及びLAN 108の両方を含むようにユーザ・サイト101を実装することができることを当業者は理解するであろう。

【0021】当業者には周知のように、SANは共有ストレージ・デバイスの高速サブネットワークであり、SANによりLAN又はWAN内のすべてのサーバからすべてのストレージ・デバイスが利用可能である。SANにさらにストレージ・デバイスを追加すると、その追加されたストレージ・デバイスも、さらに大きなネットワーク内の任意のサーバからアクセスできるようになる。

【0022】SAN 107は、たとえば、ファイバ・チャネル又は小型コンピュータ用周辺機器インタフェース (SCSI) で構成できる。当業者には周知のように、ファイバ・チャネルはコンピュータ及び大容量記憶デバイス製造業者のコンソーシアムが開発したシリアル・データ転送アーキテクチャであり、現在は、米国規格協会 (ANSI) によって標準化されている。最も有名なファイバ・チャネル標準は、Fibre Channel Arbitrated Loop (FC-AL) で、新しい大容量記憶デバイスや非常に高い帯域幅を必要とするその他の周辺機器用に設計されている。光ファイバでデバイスを接続するFC-ALは、全二重データ転送を行い、転送速度は毎秒約100メガバイト (MBps) である。

【0023】また当業者には周知のように、SCSIは種々のコンピュータシステムで周辺装置をコンピュータに接続するために使用されているパラレル・インタフェース標準である。SCSIインタフェースは、標準装備のシリアル・ポート及びパラレル・ポートに比べて高速なデータ転送速度 (最大約80MBps) を利用できる。

【0024】また当業者には周知のように、LANとは、通常比較的狭い地域に置かれるコンピュータ・ネットワークのことである。ほとんどのLANは、単一の建物又は建物群に制限されている。ほとんどのLANは、ワークステーション及びパーソナル・コンピュータを接続する。LAN内の各ノード (個々のコンピュータ)

は、個々に中央処理装置 (CPU) を備え、プログラムを実行するが、LAN上の任意の場所にあるデータ及びデバイスにアクセスすることも可能である。したがって、多くのユーザが、レーザ・プリンタなどの高価なデバイスをデータとともに共有できる。ユーザはさらに、LANを使用して、たとえば電子メールを送信したりチャット・セッションに参加したりして互いに通信することもできる。LANにはいろいろな種類のものがあるが、PCにはイーサネット (登録商標) が最も一般的である。LANは非常に高速にデータを送信することができる。しかし、LANの配線距離は制限されており、また単一のLANに接続できるコンピュータの台数にも制限がある。

【0025】図1では、LAN 108をたとえばイーサネット (登録商標) として構成している。第1のホスト104は、たとえば、ブロック (SCSI) インタフェース (I/F) を使用してローカル・ストレージ106にアクセスする。第2のホスト105は、たとえば、ファイバ (NFS又はCIFS) I/Fを使用してローカル・ストレージ106にアクセスする。第2のホスト105は、たとえば、ブロック (iSCSI) I/Fを使用してローカル・ストレージ106にアクセスすることもできる。

【0026】ローカル・ストレージ106及びサービス・プロバイダ・サイト102及び103では、ATM上でiSCSI又はファイバ・チャネル、SONET上でファイバ・チャネル、又は独自ベンダ・プロトコルを使用できる。ローカル・ストレージ106は、仮想ストレージ111、コンフィギュレーション・テーブル112 (図4も参照)、統計情報113 (図5も参照)、アクセス・ログ114 (図6も参照)、キャッシュ情報115、及びプログラム116を含む。コンポーネント112、113、114、115、及び116ではローカル・ストレージ106をキャッシュとして動作させることができる。ホスト104及び105は仮想ストレージ111にアクセスできる。仮想ストレージ111は、いくつかのボリューム117、118、119、及び120を含む。仮想ストレージ111内のボリュームの数は異なってもよい。これらのボリューム117~120の管理は、ローカル・ストレージ106が行う。ローカル・ストレージ106及びサービス・プロバイダ・サイト102及び103は、仮想ボリューム・ストレージ111の作成のため共同作業を行う。

【0027】ユーザは、ユーザ・サイト101に仮想ボリューム・ストレージ111 (ローカル・ストレージ106内) を配備することができる。仮想ボリューム・ストレージ111がユーザ・サイト101に配備されると、ホスト104及び/又はホスト105のユーザは仮想ボリューム・ストレージ111を使用することによ

り、巨大なサービス・プロバイダ・ストレージ（ボリューム121～124及び／又はボリューム125～128）にアクセスできる。仮想ボリューム・ストレージ111は通常、サービス・プロバイダ・ストレージ・ボリューム（ボリューム121～124及び／又はボリューム125～128）に比べてサイズが小さい。仮想ボリューム・ストレージ111は、たとえば以下で説明するようなブロックI/F（SCSI）、ファイルI/F（NFS、CIFS）などのいくつかのインタフェースを備える。仮想ボリューム・ストレージ111は、ユーザがWAN 110に毎回アクセスしなくてよいようにすると都合がよい。仮想ボリューム・ストレージ111ではさらに、サービス・プロバイダがユーザ向けに高速で信頼性の高いストレージ・システムを提供することもできる。

【0028】図2は、ローカル・ストレージ106の一実施形態の詳細を示すブロック図である。ローカル・ストレージ106は、マイクロプロセッサ201（たとえば、Intel CorporationやMotorola Corporation製）、ローカル（内部）バス202、SAN 107に結合されているSANインタフェース（I/F）203（図1）、LAN 108に結合されているLAN I/F 204（図1）、WAN 110に結合されているWAN I/F 205（図1）、メモリ206、メモリ207、ディスク・コントローラ208、及びローカル・ディスク209を備える。メモリ207は、コンフィギュレーション・テーブル112、統計情報113、アクセス・ログ114、キャッシュ情報115、及びプログラム116を格納する。マイクロプロセッサ201は、ローカル・ストレージ106内のすべてのリソースを制御し、プログラム116を使用してローカル・ストレージ106内のすべてのアプリケーションを実行する。

【0029】図2は、ローカル・ストレージ106内のアプリケーションを実行しているときのローカル・ディスク209のスナップショットを示している。当業者には周知のとおり、スナップショットとは、実際のストレージの全ボリューム又はボリュームの一部のコピーである。図2に示されている情報及びプログラムはローカル・ディスク209に格納されている。ローカル・ストレージ106のブート手順で、これらのデータ及びプログラムはローカル・ディスク209からメモリ207に移動される。

【0030】キャッシュ情報115は、キャッシュ・ディレクトリ210、保留データ・リスト211、及びキャッシュ・データ212で構成される。キャッシュ・ディレクトリ210は、メモリ207及びローカル・ディスク209内のキャッシュされたデータ212のディレクトリ（コンフィギュレーション情報）である。このキャッシュ・データ212は、LRU（Least Recently U

sed）管理法で管理することができる。しかし、近い将来再びアクセスする可能性がないことで、LRU管理の例外となるケースがある。このような例外の1つに順次アクセスがあり、ストレージ・アドレスが逐次的にアクセスされる。

【0031】保留データ・リスト211は、ローカル・ストレージ106に保存されている保留データのリストである。保留データ・リスト211は、保留データへのポインタ、サービス・プロバイダ・サイトへのポインタ、及び同期期間などの各保留データの情報を含む。同期期間は、データの特徴により定義される。データが重要でない場合、このデータは所定の期間、ローカル・ストレージ106内に格納することができる。この期間は、たとえば、1分、1時間、1日、1週間、又は1か月とすることができる。データが重要な場合、データが仮想ストレージ111に格納された後、コンフィギュレーション・テーブル112（図4）を使用してデータを即座に（同期して）サービス・プロバイダ・サイト（たとえば、サイト102）に格納すべきである。たとえば、同期データはボリューム00 117に格納される。図4で、ボリューム00 117は識別番号ID 00で識別される。ボリューム00 117に格納されているデータは同期データなので、データはサービス・プロバイダ・サイト102（リモート・サイト1）のボリューム10 121（リモートID 10）に即座に格納される。

【0032】この同期アクセス機能を使用すると、信頼性の高いデータ・ストレージを構築できる。ただし、ユーザ・サイト101からサービス・プロバイダ・サイト102及び／又は103にアクセスするのに時間がかかるため、同期アクセスにはかなりの時間が必要である。必要なアクセス時間は、システム100で実行される特定のオペレーションによって異なる。

【0033】ホスト104及び／又は105からアクセスされたデータはメモリ207に格納される。このデータは、キャッシュ・データ212と呼ばれる。キャッシュ・データ212は、ホスト104及び／又はホスト105のユーザによって頻繁にアクセスされる一時的データである。

【0034】一実施形態では、プログラム116は、サーバ・プログラム213、シンクロナイザ217、キャッシュ制御218、スケジューラ220、及びデバイス・ドライバ219を含む。サーバ・プログラム213では、ローカル・ストレージ106及びホスト104（及び／又はホスト105）が互いに通信できる。サーバ・プログラム213は、NFSサーバ214、CIFSサーバ215、及びSCSIサーバ216で構成される。NFSサーバ214は、従来のNFSサーバとして動作する。CIFSサーバ215、従来のCIFSサーバとして動作する。SCSIサーバ216は、SCSIター

ゲット・デバイスとして動作する。シンクロナイザ217は、保留データ・リスト211を使用して、ローカル・ストレージ106とサービス・プロバイダ・サイト102及び/又は103の同期をとる。キャッシュ制御218は、キャッシュ・ディレクトリ210を使用して、メモリ207とローカル・ディスク209内のキャッシュ・データ212を制御する。スケジューラ220は、サーバ・プログラム213、シンクロナイザ217、キャッシュ制御218、及びデバイス・ドライバ219などのすべてのプロシージャをスケジュールする。デバイス・ドライバ219は、SAN I/F 203、LAN I/F 204、WAN I/F 205、メモリ I/F 206、及びディスク・コントローラ208などのローカル・ストレージ106内のすべてのデバイスを制御する。

【0035】図3は、本発明の実施形態によるデータ・レイアウトの一例を示すブロック図である。ローカル・ストレージ106の管理は、以下でさらに詳しく説明しているように、データの特徴を使用することに基づいている。仮想ボリューム・ストレージ111では、データが静的かキャッシュされているか、同期か非同同期に応じてデータの特徴を取り扱う。以下でさらに詳しく説明しているように、仮想ボリューム・ストレージ111は、ストレージ・ボリュームだけでなく、ディレクトリ、ファイル、シリンダ、及び/又はブロック・アドレスでもデータの特徴を取り扱える。ユーザ・サイト101では、仮想ストレージ111に4種類の仮想ボリューム(ボリューム00 117、ボリューム01 118、ボリューム02 119、及びボリューム03 120)がある。これらの仮想ボリュームの特徴である、(1) 静的同期ボリューム、(2) 静的非同同期ボリューム、(3) キャッシュされた同期ボリューム、及び(4) キャッシュされた非同同期ボリュームについて次に説明する。

【0036】(1) 静的同期ボリューム(ボリューム00 117): ボリューム00 117は静的同期ボリュームである。したがって、ボリューム00 117は、第1のサービス・プロバイダ・サイト102において実際のストレージ・ボリュームであるボリューム10121と同じサイズを占有する。ホスト104及び/又は105(ユーザ・サイト101)は、ボリューム00 117からデータを読み込み、ボリューム10121からは読み込まない。アクセスされたデータがキャッシュ・データ・キャッシュ212(図2)内にある場合、ローカル・ストレージ106(図2)はローカル・ディスク209(図2)にアクセスせず、メモリ207にアクセスするだけである。ホスト104及び/又は105は、ボリューム00 117とボリューム10121にデータを書き込む。この静的非同同期ボリューム、ボリューム00 117は、高速な読み込みアクセ

スと信頼性の高いストレージ・システムを実現する。

【0037】(2) 静的非同同期ボリューム(ボリューム01 118): ボリューム01 118は静的非同同期ボリュームである。したがって、ボリューム001 118は、第1のサービス・プロバイダ・サイト102において実際のストレージ・ボリュームであるボリューム11 122と同じサイズを占有する。ホスト104及び/又は105は、ボリューム01 118からだけデータを読み込み、ボリューム11 122からはデータを読み込まない。アクセスされたデータがキャッシュ・データ212内にある場合、ローカル・ストレージ106はローカル・ディスク209にアクセスせず、メモリ207にアクセスするだけである。ホスト104及び/又は105は、ボリューム01 118にデータを書き込み、保留データ・リスト211(図2)に登録する。この保留データは、以下で説明するように、所定のスケジュールで同期処理される。

【0038】この保留データは、バックグラウンド・ジョブでローカル・ストレージ106(図2)内に圧縮し、WAN 110のトラフィックを低減することができる。この静的非同同期ボリューム、ボリューム01 118は、ボリュームの高速な読み書きを行うが、信頼度は同期ボリュームの場合と同じではないことがある。

【0039】(3) キャッシュされた同期ボリューム(ボリューム02 119): ボリューム02 119はキャッシュされた同期ボリュームである。したがって、ボリューム02 119は、第1のサービス・プロバイダ・サイト102において実際のストレージ・ボリュームであるボリューム12 123と同じサイズを占有する。ボリューム02 119はキャッシュとして動作する。したがって、ホスト104及び/又は105によってアクセスされるデータがボリューム02 119内にはない。アクセスされたデータがローカル・ディスク209(図2)内にはない場合、ローカル・ストレージ106(図2)は第1のサービス・プロバイダ・サイト102からデータを読み込んで、読み込んだデータをローカル・ディスク209に書き込む。データがローカル・ストレージ106内に入った後、ローカル・ストレージ106はメモリ207を介してローカル・ディスク209からデータをホスト104及び/又は105に送る。アクセスされたデータがキャッシュ・データ212内にある場合、ローカル・ストレージ106はローカル・ディスク209にアクセスせず、メモリ207にアクセスするだけである。ホスト104及び/又は105は、ボリューム02 119とボリューム12123にデータを書き込む。このキャッシュされた同期ボリューム、ボリューム02 119は都合よく、ローカル・ストレージ106のサイズを減らすことができる。この仮想ボリューム、ボリューム02 119ではさらに、信頼性の高いストレージ・システムを構築できる。

【0040】(4) キャッシュされた非同期ボリューム
(ボリューム03 120): ボリューム03 120
はキャッシュされた非同期ボリュームである。したがって、ボリューム03 120は、第1のサービス・プロバイダ・サイト102において実際のストレージ・ボリュームであるボリューム13 124と同じサイズを占有する。ボリューム04 120はキャッシュとして動作する。したがって、ホスト104及び/又は105によってアクセスされるデータがボリューム03 120内にない。アクセスされたデータがローカル・ディスク209内にない場合、ローカル・ストレージ106は第1のサービス・プロバイダ・サイト102からデータを読み込んで、この読み込んだデータをローカル・ディスク209に書き込む。データがすでにローカル・ストレージ106内に入っていた場合その後、ローカル・ストレージ106はメモリ207を介してローカル・ディスク209からデータをホスト104及び/又は105に送る。アクセスされたデータがデータ・キャッシュ207内にある場合、ローカル・ストレージ106はローカル・ディスク209にアクセスせず、メモリ207にアクセスするだけである。ホスト104及び/又は105は、ボリューム01 118にデータを書き込み、保留データ・リスト211 (図2) に登録する。この保留データは、あるスケジュールで同期処理される。このキャッシュされた非同期ボリューム、ボリューム03 120は都合よく、ローカル・ストレージ106のサイズを減らすことができるが、信頼性は同期ボリュームと同じにならない場合がある。

【0041】図3はさらに、サービス・プロバイダ・サイトで提供しているサービスも示している。第1のサービス・プロバイダ・サイト102のボリューム11 122は、第2のサービス・プロバイダ・サイト103内のボリューム21 126に複製される。ボリューム12 123は、第2のサービス・プロバイダ・サイト103内のボリューム22 127に複製される。ボリューム13 124は、第2のサービス・プロバイダ・サイト103内のボリューム23 128に複製される。ボリューム・データのこの複製 (リモート・コピー130) は、障害回復方法となる。ボリューム・データは、たとえば、従来のリモート・ミラーリング技術を使用して複製できる。さらに、第1のサービス・プロバイダ・サイト102では、ユーザ・サイト101の機能がなくてもボリューム12 123のデータのバックアップを実行できる (矢印132を参照)。このバックアップ・サービスにより、ユーザはデータ・バックアップの作業負担を軽減できる。

【0042】以下で説明しているように、仮想ボリューム・ストレージ111では、ホスト104及び/又は105からのアクセス活動に基づいてトレース・データを作成できる。仮想ボリューム・ストレージ111は、サ

ービス・プロバイダに対して、ユーザ・アクセス・パターンの分析のためこのトレース・データを用意することができる。ユーザ・アクセス・パターンを分析した後、サービス・プロバイダは仮想ボリューム・ストレージ111を都合よくチューニングできる。

【0043】図4は、本発明の実施形態によるコンフィギュレーション・テーブル112の一例を示すブロック図である。コンフィギュレーション・テーブル112は、パラメータ「ID」、「インタフェース」、「サブエリア」、「リモート・サイト」、「リモートID」、「ボリューム・タイプ」、「サービス」、「バックアップ」、「分析」、「サイズ」、「合計」、「使用済み」、「空き」、及び「ポリシー」で構成される。

「ID」は、ローカル・ストレージ106内の仮想ボリュームのボリューム識別を示す。たとえば、ボリューム00 117 (図3) のIDは00である。「インタフェース」は、SCSI、NFS、及びCIFSなどの仮想ボリューム・インタフェース・タイプを示す。「サブエリア」は、仮想ボリューム内のサブエリアの個数を示す。サブエリアは、ブロック・アドレス、シリンダ、ファイル、及び/又はディレクトリのあるエリアとすることができる。「ボリューム・タイプ」は、サブエリア間で異なることがある。この場合、図4の実例で、各サブエリアはボリューム・タイプが同じである。たとえば、ボリューム00 117 (図3) (IDは00) は上述のように同じ静的同期ボリューム・タイプである。「リモート・サイト」は、リモート・サイトの識別子を示す (たとえば、サービス・プロバイダ・サイト102の識別子)。これは、httpアドレスのIPアドレスでもよい。「リモートID」は、サービス・プロバイダ・サイト内のボリュームIDを示す。たとえば、サービス・プロバイダ・サイト102のボリューム10 121はリモートID 10を持つ。「ボリューム・タイプ」は、静的な又はキャッシュされた、同期又は非同期などのボリュームのタイプを示す。「サービス」は、バックアップ・サービスや分析サービスなどのサービス・プロバイダが提供するサービスを示す。サービス・プロバイダは、分析サービスを提供するときに、統計情報113 (図2) へのポインタを設定し、統計情報 (アクセス・ログ) 114を作成し、データの特徴を取得する。分析サービスが提供される場合、図4の例に示されているように「Ptr」がコンフィギュレーション・テーブル112の分析セクションで表示される。たとえば、コンフィギュレーション・テーブル112のボリューム02 119 (ID02) は、分析セクションにポインタPtrを持つ (行400を参照)。(ボリューム02 119) ID 02に対するこのポインタは、図5内の統計情報113の中のID 02を指している。ID 02は、サブエリア (ディレクトリ) 「usra」 (行500を参照) を持ち、ポインタ (「ptr」) はログPtr

セクション内にある。このポインタ (ptr) は、図6の例のボリューム02 119に対し作成されたアクセス・ログであるアクセス・ログ114を指している。

【0044】さらに、図5の例では、(ボリューム02 119の) ID 02はディレクトリ「/usrb」(行505参照)を含み、「Null」値(ポインタなし)がログPtrセクション内にある。さらに、ID 02は、ディレクトリ「/usrc」(行510を参照)を含み、Null値(ポインタなし)がログPtrセクション内にある。したがって、ディレクトリ/usrb及び/usrcにアクセス・ログを指すポインタがないため、これらのディレクトリに対し関連するアクセス・ログ114が作成されていない。

【0045】コンフィギュレーション・テーブル112で、「サイズ」は、ギガバイト単位で合計メモリ・サイズ(「合計」)を示し、また使用済みサイズ(「使用済み」)及び空きサイズ(「空き」)を示している。ボリュームがファイル・システムのローカル・ストレージ106内にない場合、合計サイズの判明している必要がある。パラメータ「ポリシー」は、同期のスケジュールを示す。たとえば、行410では、ボリューム00 117(IDは00)は、図4で「ポリシー」パラメータに「Null」値が設定されているが、このNull値は同期スケジュールがボリューム00 117に対し設定されていないことを示す。ボリューム02 119(IDは02)は、「ポリシー」パラメータで「週」値を設定しており、この値は、ボリューム02 119内のデータが週に1度サービス・プロバイダ・サイト(たとえば、サイト102)のデータと同期処理されていることを示す。ボリューム03 120(IDは03)は、「ポリシー」パラメータで「日」値を設定しており、この値は、ボリューム03 120内のデータが毎日サービス・プロバイダ・サイト(たとえば、サイト102)のデータと同期処理されていることを示す。

【0046】コンフィギュレーション・テーブル112のパラメータは、コンフィギュレーション・テーブル112内のコンフィギュレーション・データを設定することにより設定される。各ストレージ・サブシステムは関連するコンフィギュレーション・テーブルを持つ。

【0047】図5は、本発明の実施形態による統計情報113の一例を示すブロック図である。ユーザ又はサービス・プロバイダがコンフィギュレーション・テーブル112(図4)で「サービス分析」パラメータを「Ptr」に設定した場合、次のことが可能になる。ローカル・ストレージ106は、統計情報113のこのテーブルでローカル・ストレージ106のユーザ統計アクセス情報を収集する。「統計情報」は、「ID」、「サブエリア」、「IO/s」及び「MB/s」(「平均、最大読み込み」、「平均、最大書き込み」)、「ヒット率」、及び「ログPtr」で構成される。パラメータ「ID」

及び「サブエリア」(図5)は、コンフィギュレーション・テーブル112(図4)の「ID」及び「サブエリア」と同じ意味を持つ。「IO/s」は、1秒あたりの読み込み及び書き込みコマンドなど、1秒あたりのホスト/ユーザ入力又は出力活動の個数を示す。ボリューム又はボリューム・ディレクトリ内の読み込み又は書き込みアクセスは、知られている適当な方法で記録できる。

【0048】「MB/s」は1秒あたりのバイト数を示す。ローカル・ストレージ106は、ボリューム(又はボリューム内のサブエリア)での読み込み及び書き込みの発生数を収集し、この発生数の平均(Ave)をとる。ローカル・ストレージ106は、さらに、ボリューム又はボリューム内のサブエリアでの読み込み及び書き込みの発生数の最大発生回数(Max)も追跡する。

【0049】「ヒット率」は、ローカル・ストレージのヒット率を示す。このヒット率は、たとえば、ローカル・ディスク209(図2)だけのものである。ヒット率は、次の式(1)で定義される。

(1) ヒット率 = 100% (#ローカル・ディスク読み込み / #全読み込み)

ただし、「#ローカル・ディスク読み込み」は、ホスト別のローカル・ディスク209内の読み込み回数であり、「#全読み込み」は、ホスト別の合計読み込み回数である。

【0050】統計情報113は、サービス・プロバイダ・サイト102及び/又は103に定期的に送られる。サービス・プロバイダでは、データの特徴を分析し、この分析結果からサービス・プロバイダはユーザ・サイト101のホストのユーザに適切な解決策を提示することができる。

【0051】上述のように、ログPtrは、アクセス・ログ114へのポインタを示している。サービス・プロバイダは、ユーザ・アクセス・パターンの詳細を知りたい場合、このポインタを作成されたアクセス・ログ114に設定する。

【0052】図6は、本発明の実施形態によるアクセス・ログ114の一例を示すブロック図である。サービス・プロバイダは、ユーザのアクセス・パターンの詳細を分析するときに、アクセス・ログ114を作成し、ログPtr(図5の統計情報113内にある)のポインタをこのアクセス・ログ114に設定する。ローカル・ストレージ106は、ユーザ別に各アクセスのアクセス・ログ114を収集する。アクセス・ログ114は、パラメータ「日付」、「時刻」、「コマンド」、「ファイルID」、「アドレス」、及び「サイズ」で構成される。

「日付」は、ユーザがアクセスした日付を示す。「時刻」は、アクセスの時刻を示す。「コマンド」は、アクセスのコマンド・タイプ(たとえば、読み込み又は書き込み)を示す。「ファイルID」は、このコマンドでアクセスしたファイルの識別を示す。アクセス・ログがS

CSI用であれば、「ファイルID」はNullとすることもできる。「アドレス」及び「サイズ」は、アクセス・アドレス及びサイズを示す。

【0053】読み込みプロセス

読み込みプロセスでは、キャッシュされたデバイス（ローカル・ストレージ106）はキャッシュとして動作する。ホスト104及び／又は105は、ローカル・ストレージ106にアクセスするが、これは、ローカル・ストレージ106の空き容量が大きいためである。まず、ホスト（たとえば、ホスト104又はホスト105）は、LAN 108（図1）を介してNFS、CIFS、又はiSCSIプロトコルの読み込みコマンドを発行するか、又はSAN 107を介してSCSIプロトコルの読み込みコマンドを発行する。ローカル・ストレージ106は、LAN I/F 204又はSAN I/F 203（図2）を介して読み込みコマンドを受信する。デバイス・ドライバ219（図2）は、ホストからのこのコマンドを処理し、このコマンドをスケジューラ220に入れる。次にスケジューラ220は読み込みコマンドを分析して、読み込みコマンドをサーバ・プログラム213内において適切なサーバ（NFSサーバ214、CIFSサーバ215、又はSCSIサーバ216）に入れて読み込みコマンドを処理することができる。各サーバ・プログラム213は、読み込みコマンドで要求されたデータがキャッシュ・データ212（図2）内にあるかどうかを調べる。すべてのデータ（読み込みコマンドで要求された）がキャッシュ・データ212内にある場合は、ローカル・ストレージ106は内部バス202とネットワーク・インタフェース（SAN I/F 203又はLAN I/F 204）を介してデータを要求側ホスト（ホスト104又は105）に返す。他方、データ（読み込みコマンドで要求された）の全部又は一部がキャッシュ・データ212内にない場合、要求されたデータはローカル・ディスク209から、又はサービス・プロバイダ・サイト102又は103からキャッシュ・データ212に移動すべきである。

【0054】図7は、本発明の実施形態によるこの読み込みプロセスの流れ図である。ローカル・ストレージ106がホスト（たとえば、ホスト104又は105）から読み込みコマンドを受け取った後、ローカル・ストレージ106はデータのボリューム・タイプを調べる（ステップ701）。スケジューラ220は、ボリュームIDをチェックすることによりコンフィギュレーション・テーブル112内のボリューム・タイプ（たとえば、静的タイプ）を調べる。特に、スケジューラ220（図2）は、読み込みコマンドを分析して、データのボリューム・タイプを判別し、また読み込みコマンドを処理する（取り扱う）サーバ・プログラム（サーバ214、215、又は216）を決定する。スケジューラ220は、コンフィギュレーション・テーブル112（図4）

を見て読み込みコマンドを処理するためのボリューム・タイプと適切なサーバ・プログラムを調べる。ボリューム・タイプが静的ボリュームの場合、ローカル・ストレージ106はステップ702、703、及び704を飛ばし、以下で説明するステップ705に進む。ボリューム・タイプがキャッシュされたボリュームの場合、ローカル・ストレージ106はデータ（読み込みコマンドで要求した）がキャッシュ・データ212内に格納されているかないかを調べる（ステップ702）。

【0055】ローカル・ストレージ106は、キャッシュ・ディレトリ210を調べる（ステップ703）。特に、サーバ・プログラム213の適切なサーバ（214、215、又は216）は、読み込みコマンドで要求したデータがローカル・ディスク209内にあるかどうかを調べる。すべてのデータ（読み込みコマンドで要求した）がローカル・ディスク209内に格納されている場合、ローカル・ストレージ106はステップ704を飛ばして、以下で説明するステップ705に進む。データ（読み込みコマンドで要求した）がローカル・ディスク209内にない場合、ローカル・ストレージ106はステップ704を実行する。

【0056】ステップ704で、ローカル・ストレージ106はローカル・ディスク209内でデータ領域を割り当て、サービス・プロバイダ・サイト（たとえば、サービス・プロバイダ・サイト102）からデータ（読み込みコマンドで要求した）を読み込む。特に、キャッシュ制御218では、コンフィギュレーション・テーブル112（図4）を使用してデータが取り出されるサービス・プロバイダ・サイトを決定する。キャッシュ制御218では、サービス・プロバイダ・サイトからデータを読み込み、データをローカル・ディスク209に格納する。

【0057】データをローカル・ディスク209に格納した後、ローカル・ストレージ106はローカル・ディスク209からキャッシュ・データ212へデータを移動する（書き込む）（ステップ705）。特に、キャッシュ制御218は、ローカル・ディスク209からキャッシュ・データ212にデータを移動する。続いてローカル・ストレージ106は、統計情報113（図5）を更新する（ステップ706）。特に、ローカル・ストレージ106のスケジューラ220は、統計情報113を更新する。統計情報113内のログ・ポイント（ログポイント）（図5）が設定されている場合、ローカル・ストレージ106はこの読み込みコマンドのログ・データをアクセス・ログ114（図6）に追加する。スケジューラ220又はデバイス・ドライバ219は、統計情報113とアクセス・ログ114を作成する。通常、統計情報113とアクセス・ログ114の作成においてはスケジューラ220が用いられる。

【0058】その後、データ（読み込みコマンドで要求

した)はキャッシュ・データ212から、読み込みコマンドを送ったホスト(たとえばホスト104又は105)に移動される。特に、キャッシュ制御218はキャッシュ・データ212からデータを適切なサーバ・プログラム(214、215、又は216)に移動し、適切なサーバ・プログラム(214、214、又は216)がデータを、読み込みコマンドを送ったホストに返す。

【0059】書き込みプロセス

書き込みプロセスでは、キャッシュされたデバイス(ローカル・ストレージ106)はキャッシュとして動作する。ホスト104及び/又は105は、ローカル・ストレージ106にアクセスするが、これは、ローカル・ストレージ106の空き容量が大きいかからである。まず、ホスト(たとえば、ホスト104又はホスト105)は、LAN 108を介してNFS、CIFS、又はSCSIプロトコルの書き込みコマンドを発行するか、又はSAN 107を介してSCSIプロトコルの書き込みコマンドを発行する。ローカル・ストレージ106は、LAN I/F 204又はSAN I/F 203を介して書き込みコマンドを受信する。デバイス・ドライバ219(図2)は、ホストからのこの読み込みコマンドを処理し、このコマンドをスケジューラ220に入れる。次にスケジューラ220は書き込みコマンドを分析して、書き込みコマンドをサーバ・プログラム213内において適切なサーバ(NFSサーバ214、CIFSサーバ215、又はSCSIサーバ216)に入れる。各サーバ・プログラム213は、書き込みコマンドのデータがキャッシュ・データ212(図2)内にあるかどうかを調べる。すべてのデータ(書き込みコマンドの)がキャッシュ・データ212内にある場合は、ローカル・ストレージ106は内部バス202とネットワーク・インタフェース(SAN I/F 203又はLAN I/F 204)を介してホストからデータを受信する。ローカル・ストレージ106は、ローカル・ディスク209及びキャッシュ212内の割り当てられた空き領域にデータを上書きする。他方、(書き込みコマンドの)データの全部又は一部がキャッシュ・データ212内がない場合、ローカル・ストレージ106は、キャッシュ・データ212内がないデータの残りについてある領域(ローカル・ディスク209及びキャッシュ・データ212内の)を割り当てる。すべてのデータ領域の割り当てが済んだら、ローカル・ストレージ106はその割り当てられた領域内にデータを格納する。すべてのデータをキャッシュ・データ212内に格納する場合、ローカル・ストレージ106はローカル・ディスク209にデータを格納し、このデータをプロバイダ・サイト(たとえば、サービス・プロバイダ・サイト102)に送信する。

【0060】図8は、本発明の実施形態による書き込みプロセスの流れ図である。ローカル・ストレージ106

がホストから書き込みコマンドを受け取った後、ローカル・ストレージ106は書き込みコマンドのデータのボリューム・タイプを調べる(ステップ801)。特に、スケジューラ220(図2)は、書き込みコマンドを分析して、データのボリューム・タイプを判別し、また読み込みコマンドを処理する(取り扱う)サーバ・プログラム(サーバ214、215、又は216)を決定する。スケジューラ220は、コンフィギュレーション・テーブル112(図4)を見て書き込みコマンドを処理するためのボリューム・タイプと適切なサーバ・プログラムを調べる。ボリューム・タイプが静的ボリュームの場合、ローカル・ストレージ106はステップ802、803、及び804を飛ばし、以下で説明するステップ805に進む。ボリューム・タイプがキャッシュされたボリュームの場合、ローカル・ストレージ106はデータ(書き込みコマンドの)がキャッシュ・データ212(図2)内に格納されているかいないかを調べる。ローカル・ストレージ106は、キャッシュ・ディレクトリ210(図2)を調べる(ステップ802)。

【0061】特に、サーバ・プログラム213の適切なサーバ(214、215、又は216)は、書き込みコマンドのデータがローカル・ディスク209内にあるかどうかを調べる。ステップ803で、すべてのデータ(書き込みコマンドで要求した)がローカル・ディスク209内で割り当てられている場合、ローカル・ストレージ106はステップ804を飛ばして、ステップ805に進む。ステップ803で、データがローカル・ディスク209内にない場合、ローカル・ストレージ106はステップ804を実行する。

【0062】(書き込みコマンドの)データがローカル・ディスク209内で割り当てられていない場合、ローカル・ストレージ106は書き込みデータについてローカル・ディスク209とキャッシュ・データ212の両方にデータ領域を割り当てる(ステップ804)。特に、キャッシュ制御218は、ローカル・ディスク209及びキャッシュ・データ212の両方でデータを割り当てる。

【0063】ローカル・ストレージ106は、キャッシュ・データ212を介して書き込みデータをローカル・ディスク209に書き込む。特に、適切なサーバ・プログラム(NFSサーバ214、CIFSサーバ215、又はSCSIサーバ216)はデータをローカル・ディスク209に書き込む。

【0064】ローカル・ストレージ106は、コンフィギュレーション・テーブル112(図4)を使用してデータを書き込む際のボリューム・タイプを調べる(ステップ806)。特に、サーバ・プログラム(NFSサーバ214、CIFSサーバ215、又はSCSIサーバ216)はコンフィギュレーション・テーブル112を使用してボリューム・タイプを調べる。ボリューム・タ

イブが同期領域の場合、ローカル・ストレージ106は即座にサービス・プロバイダ・サイト（たとえば、サイト102）とデータの同期をとり（書き込み）（ステップ808）、以下で説明するステップ809に進む。シンクロナイザ217（図2）はサービス・プロバイダ・サイトとの（ユーザ・サイト内の）データの同期処理を実行する。ローカル・ストレージ106内のキャッシュ制御218は、コンフィギュレーション・テーブル112を使用してサービス・プロバイダ・サイトに格納すべきデータを認識する。

【0065】ステップ806で、ボリューム・タイプが非同期領域の場合、ローカル・ストレージ106は保留データ・リスト211（図2）を更新する（ステップ807）。特に、サーバ・プログラム（サーバ214、215、又は216）が保留データ・リストを更新する。

【0066】続いてローカル・ストレージ106は、統計情報113（図5）を更新する（ステップ809）。特に、スケジューラ220は、統計情報113を更新する。統計情報113内のログ・ポインタ（ログポインタ）（図5）が設定されている場合、ローカル・ストレージ106はこの書き込みコマンドのログ・データをアクセス・ログ114に追加する。特に、スケジューラ220は書き込みコマンドのログ・データをアクセス・ログ114に追加する。これで書き込みプロセスは終了する。

【0067】同期プロセス

非同期ボリューム書き込み（非同期領域へのデータ書き込み）の場合、データは、ユーザ又はサービス・プロバイダによって定められたスケジュールに従ってサービス・プロバイダ・サイト（たとえば、サイト102）に送信される。このスケジュールは、たとえば、図4のコンフィギュレーション・テーブル112内の「ポリシー」エントリ内の値により定められる。スケジューラ220（図2）は、シンクロナイザ217（図2）を定期的に実行することで、サービス・プロバイダ・サイトとデータの同期をとることができる。この期間は、システム100設定に応じて、たとえば、約1.0ミリ秒又は10.0ミリ秒に設定できる。図9は、本発明の実施形態による同期プロセスの流れ図である。同期プロセスは、シンクロナイザ217（図2）によって実行できる。

【0068】シンクロナイザ217は、保留データ・リスト211（図2）内のヘッド・データを選択する（ステップ901）。最初に、シンクロナイザ217は、保留データ・リスト211を調べる（ステップ902）。保留データ・リスト211にデータがなければ、シンクロナイザ217は同期プロセスを終了する。同期すべきデータが1つ又は複数（保留データ・リスト211内）にある場合、シンクロナイザ217は以下で説明するようにステップ902〜905を実行する。

【0069】シンクロナイザ217は、保留データ・リスト211の情報を調べる（ステップ903）。保留デ

ータ・リスト211内のこの保留データの同期をとる必要があれば、シンクロナイザ217は保留データをサービス・プロバイダ・サイト（たとえば、サイト102）に送り、保留データとサービス・プロバイダ・サイトとの同期をとる（ステップ904）。ローカル・ストレージ106内のキャッシュ制御218は、コンフィギュレーション・テーブル112（図4）内の「リモートID」及び「リモート・サイト」の値に基づいてサービス・プロバイダ・サイト（たとえば、サイト102）に格納すべきデータを認識する。

【0070】シンクロナイザ217は、保留データ・リスト211内の次のデータを選択する（ステップ905）。シンクロナイザ217は、保留データ・リスト211内で選択するデータがなくなるまでステップ902〜905を繰り返す。保留データ・リスト211に選択するデータがなくなったら、図9の方法は終了する。コンピュータが上述のいずれかの方法を実行できるようにコンピュータで読み取り可能な媒体に格納できるプログラム又はコードを実装することも本発明の範囲である。

【0071】そこで、本発明について特定の実施形態を引用しながら説明してきたが、前記開示においては修正、さまざまな変更、及び置き換えを許容することが意図されており、規定されているとおり本発明の範囲を逸脱することなく他の特徴の対応する使用がなくても場合によっては本発明のいくつかの特徴を採用することは認められるであろう。

【図面の簡単な説明】

【図1】本発明の実施形態によるシステムのブロック図。

【図2】図1に示されているローカル・ストレージの一実施形態の詳細を示すブロック図。

【図3】本発明の実施形態によるデータ・レイアウトの一例を示すブロック図。

【図4】本発明の実施形態によるコンフィギュレーション・テーブルの一例を示すブロック図。

【図5】本発明の実施形態による統計情報の一例のブロック図。

【図6】本発明の実施形態によるアクセス・ログの一例のブロック図。

【図7】本発明の実施形態による読み込みプロセスの流れ図。

【図8】本発明の実施形態による書き込みプロセスの流れ図。

【図9】本発明の実施形態による同期プロセスの流れ図。

【符号の説明】

101…ユーザ・サイト、102、103…サービス・プロバイダ・サイト、104、105…ホスト、106…ローカル・ストレージ、107…SAN、108…LAN、111…仮想ストレージ

【 図 1 】

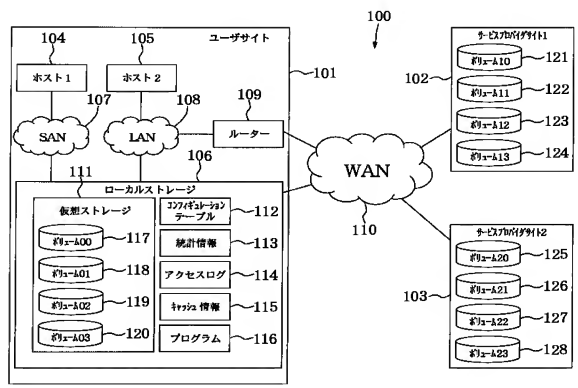


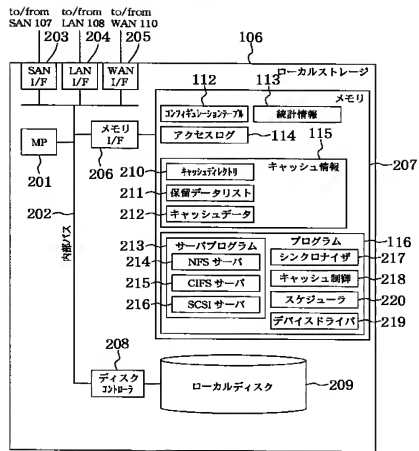
図 1

【 図 2 】

【 図 6 】

図 2

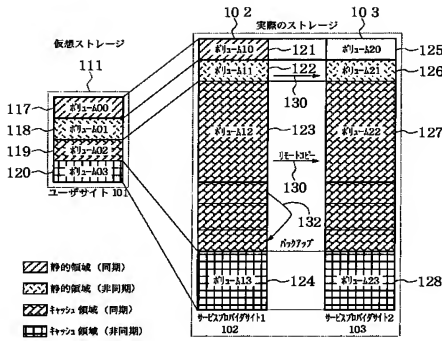
図 6



アクセスログ 114

日付	時刻	コマンド	ファイルID	アドレス	サイズ
2000.8.1	0:0:0	読み込み	0	0x0000	0x0001
2000.8.1	0:0:0	書き込み	0	0x0000	0x0001
2000.8.1	0:0:1	読み込み	10	0x0100	0x0100
2000.8.1	0:0:1	読み込み	10	0x0200	0x0100
2000.8.1	0:0:2	読み込み	10	0x0200	0x0100
2000.8.1	0:0:2	読み込み	10	0x0300	0x0100
2000.8.1	0:0:2	読み込み	10	0x0400	0x0100
...					

【図3】



【例4】

コンフィギュレーションテーブル

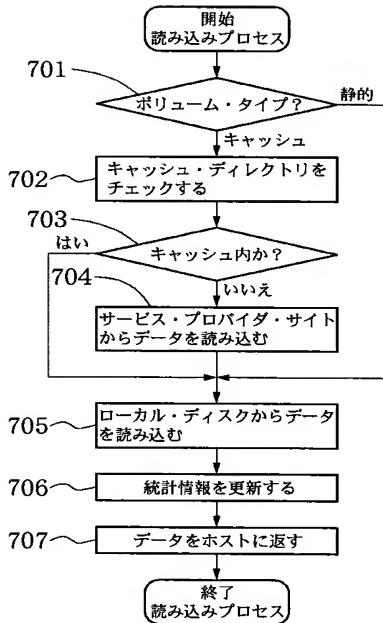
ID	インタフェース	サブエリア	リモートサイト	リモートID	ポリシータイプ	サービス		サイズ (GB)			ポリシー
						分析	バックアップ	合計	使用済み	空き	
00	SCSI	None	1	10	静的	リモートコピー	Null	20			Null
01	SCSI	None	1	11	静的 (非復旧)	None	Ptr	20			Null
02	NFS	3	1	12	キャッシュ (非同期)	リモートコピー	Ptr	100	50	50	遷
03	CIFS	None	1	13	キャッシュ	None	Ptr	50	30	20	日

【图5】

[illegible]

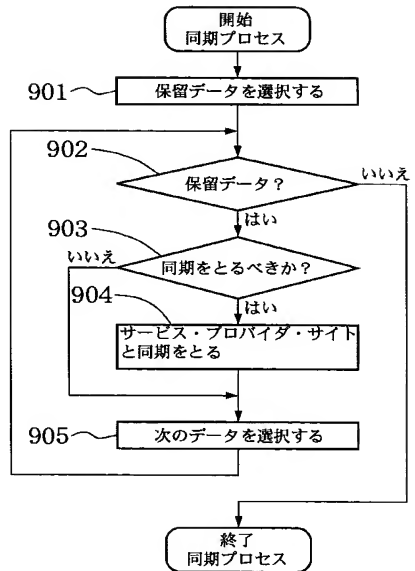
【図7】

図 7



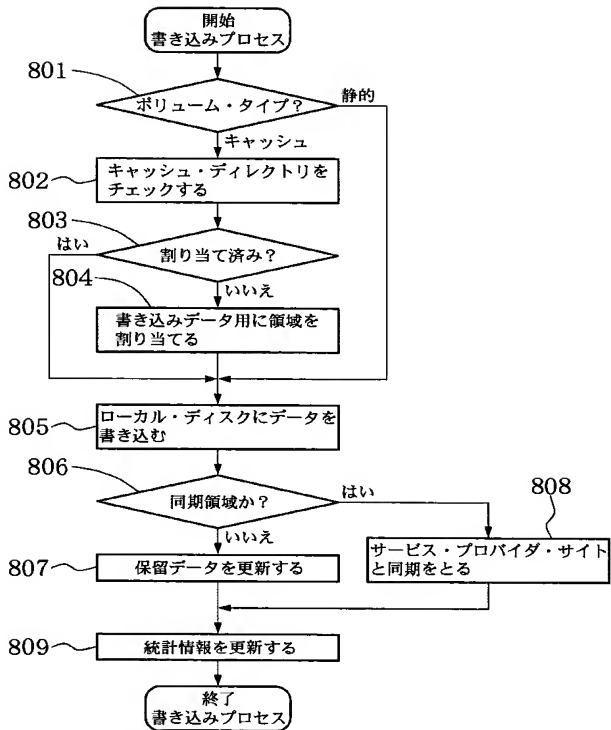
【図9】

図 9



【図8】

図 8



フロントページの続き

Fターム(参考) 5B065 CE21 EA35 EK05
5B082 FA11 FA12 GB06 HA02

【外国語明細書】

1 Title of Invention

VIRTUAL VOLUME STORAGE

2 Claims

1. A system for providing a data storage service, the system comprising
:

a service provider site configured to provide a data storage service
; and

a user site coupled by a wide area network (WAN) to the service provider site, the user site comprising a local storage having a virtual storage, the virtual storage having a synchronous volume and an asynchronous volume, the local storage configured to immediately transmit to the service provider site data that is written in the synchronous volume, to transmit at a predetermined schedule to the service provider site data that is written in the asynchronous volume, and to read data from the service provider site if the data is not stored in the local storage.

2. The system of claim 1 wherein the local storage is configured to record an access activity to the local storage by the host.

3. The system of claim 1 wherein the local storage records an access activity to a volume in the virtual storage.

4. The system of claim 1 wherein the local storage records an access activity in a sub area of a volume in the virtual storage.

5. The system of claim 2 wherein the access activity is provided to the service provider site for analysis.

6. A method of providing a data storage service, the method comprising:
providing a local storage having a virtual storage, the virtual storage comprising a synchronous volume and an asynchronous volume; and
selectively performing at least one of the following:

immediately transmitting to a service provider site data that is written in the synchronous volume of the virtual storage;

transmitting at a predetermined schedule to the service provider site data that is written in the asynchronous volume the virtual storage; and

reading data from the service provider site if the data is not stored in the local storage.

7. The method of claim 6, further comprising:

recording an access activity to the local storage.

8. The method of claim 7, further comprising:

providing a record of the access activity to the service provider site.

9. An electronically-readable medium storing a program for permitting a computer to perform a method of providing a data storage service, the method comprising:

selectively performing at least one of the following:

immediately transmitting to a service provider site data that is written in a synchronous volume of a virtual storage;

transmitting at a predetermined schedule to the service provider site data that is written in an asynchronous volume the virtual storage; and

reading data from the service provider site if the data is not stored in a local storage that includes the virtual storage.

10. A program code embedded on a carrier wave for causing a computer to perform a method of providing a data storage service, the method comprising:

selectively performing at least one of the following:

immediately transmitting to a service provider site data that is written in a synchronous volume of a virtual storage;

transmitting at a predetermined schedule to the service provider site data that is written in an asynchronous volume the virtual storage; and

reading data from the service provider site if the data is not stored in a local storage that includes the virtual storage.

11. An apparatus for processing data in a data storage system, the apparatus comprising:

a user site comprising a local storage having a virtual storage, the virtual storage including a synchronous area and an asynchronous area, the local storage configured to immediately transmit to a service provider site data that is written to the synchronous area, to transmit at a predetermined schedule to the service provider site data that is written to the asynchronous area, and to read data from the service provider site if the data is not stored in the local storage.

12. A method of processing data in a data storage system, the method comprising:

in response to a command signal to a local storage, selectively performing at least one of the following:

immediately transmitting to a service provider site data that is written to a synchronous area in a virtual storage;

transmitting at a predetermined schedule to the service provider site data that is written to an asynchronous area in the virtual storage;

reading data from the service provider site if the data is not stored in the local storage.

13. The method of reading data in a data storage system, the method comprising:

receiving a read command from a host to a local storage;

determining the volume type of the data that is requested by the read

d command;

if the volume type is static, then reading the data from a local disk in the local storage;

if the volume type is cached, then checking if the data is stored in the local disk;

if the data is stored in the local disk, then reading the data from the local disk;

if the data is not stored in the local disk, then reading the data from a remote service provider site, storing the data in the local disk, and reading the data from the local disk;

after reading the data from the local disk, updating statistical information relating to the read command; and

returning the data to the host.

14. An apparatus for processing data in a data storage system, the apparatus comprising:

a service provider site configured to immediately receive data that is written to a synchronous area in a virtual storage, to receive at a predetermined schedule data that is written to an asynchronous area in the virtual storage, and to transmit data to a local storage that includes the virtual storage if the data is not stored in the local storage.

15. A method of processing data in a data storage system, the method comprising:

selectively performing one of the following:

immediately receiving data that is written to a synchronous area in a virtual storage;

receiving at a predetermined schedule data that is written to an asynchronous area in the virtual storage; and

transmitting data to a local storage that includes the virtual storage if the data is not stored in the local storage.

16. A method of writing data in a data storage system, the method comprising:

receiving a write command from a host to a local storage;

determining the volume type of the data of the write command;

if the volume type is static, then writing the data to a local disk in the local storage;

if the volume type is cached, then checking if the data is stored in the local disk;

if the data is stored in the local disk, then writing the data to the local disk from a cache;

if the data is not stored in the local disk, then allocating a data area in the local disk and in the cache and then writing the data to the local disk from the cache;

checking if the volume type of the data is synchronous;

if the volume type is synchronous, then immediately synchronizing the data with a remote service provider site and then updating statistical information relating to the write command; and

if the volume type is not synchronous, then synchronizing the data with a remote service provider site based on a predetermined schedule and then updating statistical information relating to the write command.

17. A system for providing a data storage service, the system comprising:

means for selectively performing at least one of the following:

immediately transmitting to a service provider site data that is written in a synchronous volume of a virtual storage;

transmitting at a predetermined schedule to the service provider site data that is written in an asynchronous volume of the virtual storage; and

reading data from the service provider site if the data is not s

tored in a local storage that includes the virtual storage.

3 Detailed Description of Invention

FIELD OF THE INVENTION

The present invention relates to data storage systems, and relates more particularly to a system and method for providing a virtual volume of data storage in a user site.

BACKGROUND OF THE INVENTION

Managing a large data storage system is very difficult. Typically, a data storage service provider provides a storage-related service such as providing storage volumes, data backup, and the like. From a user site, if a user connects to this storage via the storage service provider's directory, then the user must use a long distance connection to connect to this storage. This process causes more delays than if a user is connecting to a local storage.

The World Wide Web (WWW) effectively acts as a storage system in a wide area. The user may deploy a proxy server in the user site. As known to those skilled in the art, a proxy server is a server that sits between a client application (such as a web browser) and a remote server. The proxy server provides a cache of items that are available on the remote servers. The proxy server intercepts all requests that are made to the remote server so that the proxy server can determine if it can instead fulfill the request. If the proxy server is unable to fulfill the request, then the proxy server will forward the request to the remote server for processing. A proxy server can just only read cache data from the remote server or site, and it does not support a write procedure to the remote server or site. Additionally, the proxy server can not provide a service that depends on the data feature. As a result, the proxy serv

er causes an ineffective usage of resources, and the use of a proxy server also causes more delays than the use of a local storage.

U.S. patent 5,155,835 discloses a multilevel, hierarchical, dynamically mapped data storage subsystem. This patent reference does not disclose storage systems between wide area networks (WANs) and ignores data features such as access interface type (block or file, and the like), purpose (user data or system file or backup file and the like). Thus, the subsystem disclosed in this reference is inefficient.

The Internet protocol small computer system interface (iSCSI) uses the Internet Protocol (IP) networking infrastructure to quickly transport large amounts of block storage (SCSI) data over existing local area and/or wide area networks. With the potential to support all major networking protocols, iSCSI (IP SAN) can unify network architecture across an entire enterprise, thereby reducing the overall network cost and complexity. To ensure reliability, iSCSI can use known network management tools and utilities that have been developed for IP networks. The iSCSI protocol is discussed, for example, at the website, <http://www.ece.cmu.edu/ips/index.html> in the IP Storage section, <http://www.ece.cmu.edu/ips/docs/docs.html>.

The working group of Internet Engineering Task Force (IETF) at <http://www.ietf.org> provides a network file system (NFS) version 3 (rfc1813 of IETF) and a common Internet File system (CIFS) protocol (<http://www.cifs.org>).

The NFS is an open operating system that allows all network users to access shared files that are stored in different types of computers. NFS provides access to shared files through an interface called Virtual File System (VFS) which runs on top of the Transmission Control Protocol/Internet Protocol (TCP/IP). With NFS, computers connected to a network can operate as clients while accessing remote files and as servers while

providing remote users access to local shared files.

The CIFS protocol defines a standard for remote file access using millions of computers at a time. With CIFS, users with different platforms and computers can share files without having to install new software. CIFS runs over TCP/IP, but uses the Server Message Block (SMB) protocol found in Microsoft Windows for file and printer access. Therefore, CIFS will allow all applications (including Web browsers) to open and share files across the Internet.

The websites <http://www.cnt.com> and <http://www.san.com> describe wide area network (WAN) connections to a storage or storage area networks (SANs). The WAN connection may be an asynchronous transfer mode (ATM), synchronous optical network (SONET), and the like.

The above references do not disclose methods for providing hierarchical management techniques between storage systems and host devices.

There is a need for a system and method that will overcome the above-mentioned deficiencies of conventional methods and systems. There is also a need for a system and method that will permit an access service provider to have a reliable storage system and that will permit a user to quickly access the storage system. There is also a need for a system and method that will permit an access service provider to be able to tune a local storage system based upon user statistic data and user log data.

SUMMARY

The present invention may advantageously provide a reliable virtual local storage in a user site. The present invention may also advantageously provide a virtual local storage that permits faster access for a user. The present invention may also advantageously permit a service provider to tune up the local storage by using user statistics data and user log data that are tracked by the local storage, and the service provide

r may then charge the user for these tuning services.

In one embodiment, the present invention permits a service provider to provide a data storage service via a wide area network. A virtual volume storage is deployed at the user site. The virtual volume storage allows a host(s) at the user site to access the service provider storage located at the service provider site. The virtual volume storage advantageously permits the user to avoid having to access each time the wide area network coupled between the user site and the service provider site.

In one embodiment, the present invention broadly provides a system for providing a data storage service, comprising: a service provider site configured to provide a data storage service; and a user site coupled by a wide area network (WAN) to the service provider site, the user site comprising a local storage having a virtual storage, the virtual storage having a synchronous volume and an asynchronous volume, the local storage configured to immediately transmit to the service provider site data that is written in the synchronous volume, to transmit at a predetermined schedule to the service provider site data that is written in the asynchronous volume, and to read data from the service provider site if the data is not stored in the local storage.

In another embodiment, the present invention provides a method of managing data in a local storage based on the data feature. The virtual volume storage in the local storage determines if the data is static or cached, and synchronous or asynchronous. The virtual volume storage can also manage data not only by storage volume, but also by directory, file, cylinder, and/or block address.

In another embodiment, the present invention provides a method of tracking the access activities of a user in the user site. The user access activities are recorded to permit the service provider to analyze patterns in the user access activities. Based on this analysis, the service

provider can tune the virtual volume storage in the user site.

DETAILED DESCRIPTION OF THE PREFERRED EMBODIMENTS

The following description is provided to enable any person skilled in the art to make and use the present invention, and is provided in the context of a particular application and its requirements. Various modifications to the embodiments will be readily apparent to those skilled in the art, and the generic principles defined herein may be applied to other embodiments and applications without departing from the spirit and scope of the present invention. Thus, the present invention is not intended to be limited to the embodiments shown, but is to be accorded the widest scope consistent with the principles, features, and teachings disclosed herein.

In one embodiment in accordance with the present invention, a system 100 is configured by a user site 101, and by a first service provider site 102 and a second service provider site 103. The details of the user site 101 is described below. The number of service provider sites in the system 100 may vary. The first service provider site 102 includes storage volumes 121, 122, 123, and 124, while the second service provider site 103 includes storage volumes 125, 126, 127, and 128. The number of storage volumes in a service provider site may vary. The service provider sites 102 and 103 may be located in, for example, two different safe buildings. The user site 101, first service provider site 102, and second service provider site 103 are each connected to a Wide Area Network (WAN) 110. The service providers 102 and 103 provide data storage services via the WAN 110 to a user at the user site 101.

As known to those skilled in the art, a WAN is a computer network that typically spans a relatively large geographical area. A WAN may also

include local area networks (LANs). Computers connected to a WAN are often connected through public networks, such as the telephone system. They can also be connected through leased lines or satellites. The largest WAN in existence is the Internet.

The WAN 110 permits communication between each site (user site 101 and service provider sites 102 and 103). Each site 101, 102, and 103 may be separated from each other by long distances. The WAN 110 is typically configured by asynchronous transfer mode (ATM), synchronous optical network (SONET), Dense Wavelength Division Multiplexing (DWDM), or Internet Protocol (IP) network.

At the user site 101, a first host 104 is connected to a local storage 106 via a storage area network (SAN) 107, while a second host 105 is connected to the local storage 106 via a local area network (LAN) 108. The hosts 104 and 105 access the local storage 106. The hosts 104 and 105 are, for example, servers. The number of hosts in the user site 101 may vary. The LAN 108 is connected to the WAN 110 via, for example, a router 109. From the teachings of the present invention herein, it is understood by those skilled in the art that the user site 101 may be implemented to include only one of the SAN 107 or LAN 108, or both the SAN 107 and LAN 108.

As known to those skilled in the art, a SAN is a high-speed sub-network of shared storage devices wherein the SAN makes all storage devices available to all servers in a LAN or WAN. As more storage devices are added to a SAN, these added storage devices will also be accessible from any server in the larger network.

The SAN 107 may be configured with, for example, fibre channel or Small Computer System Interface (SCSI). As known to those skilled in the art, a fibre channel is a serial data transfer architecture developed by a consortium of computer and mass storage device manufacturers and now

being standardized by the American National Standards Institute (ANSI).

The most prominent fibre channel standard is the Fibre Channel Arbitrated Loop (FC-AL) which is designed for new mass storage devices and other peripheral devices that require a very high bandwidth. Using optical fibers to connect the devices, FC-AL supports full-duplex data transfer rates of approximately 100 megabytes per second (MBps).

As also known to those skilled in the art, SCSI is a parallel interface standard used by Apple Macintosh computers, personal computers (PCs), and many UNIX systems for attaching peripheral devices to computers. SCSI interfaces provide for faster data transmission rates (up to about 80 MBps) than standard serial and parallel ports.

As also known to those skilled in the art, a LAN is a computer network that typically spans a relatively small area. Most LANs are confined to a single building or group of buildings. Most LANs connect workstations and personal computers. Each node (individual computer) in a LAN has its own central processing unit (CPU) with which it executes programs, but it is also able to access data and devices anywhere on the LAN. Thus, many users can share expensive devices, such as laser printers, as well as data. Users can also use the LAN to communicate with each other, by, for example, sending e-mail or engaging in chat sessions. There are many different types of LANs, with Ethernet being the most common for PCs. LANs are capable of transmitting data at very fast rates, much faster than the data transmitted over a telephone line. However, the distances over LANs are limited, and there is also a limit on the number of computers that can be attached to a single LAN.

In Figure 1, the LAN 108 is configured as, for example, an Ethernet.

The first host 104 accesses the local storage 106 by use of, for example, a block (SCSI) interface (I/F). The second host 105 accesses the local storage 106 by, for example, a file (NFS or CIFS) I/F. The second

host 105 may also access the local storage 106 by, for example, a block (iSCSI) I/F.

The local storage 106 and service provider sites 102 and 103 may use iSCSI or fiber channel over ATM, fiber channel over SONET, or a unique vendor protocol.

The local storage 106 includes a virtual storage 111, configuration table 112 (see also Figure 4), statistical information 113 (see also Figure 5), access log 114 (see also Figure 6), cache information 115, and programs 116. The components 112, 113, 114, 115, and 116 permit the local storage 106 to act like a cache. The hosts 104 and 105 can access the virtual storage 111. The virtual storage 111 includes some volumes 117, 118, 119 and 120. The number of volumes in the virtual storage 111 may vary. These volumes 117-120 are managed by the local storage 106. The local storage 106 and service provider sites 102 and 103 collaborate for creation of the virtual volume storage 111.

The user can deploy the virtual volume storage 111 (in local storage 106) at the user site 101. The virtual volume storage 111 allows the users of host 104 and/or host 105 to access the huge service provider storage (volumes 121-124 and/or volumes 125-128) as the virtual volume storage 111 is deployed at the user site 101. The virtual volume storage 111 typically has a smaller size than the service provider storage volumes (volumes 121-124 and/or volumes 125-128). The virtual volume storage 111 has several interfaces such as, for example, a block I/F (SCSI), file I/F (NFS, CIFS) as discussed below. The virtual volume storage 111 advantageously permits the user to avoid in having to access the WAN 110 every time. The virtual volume storage 111 also enables the service provider to provide a faster and more reliable storage system for the user.

Figure 2 is a block diagram showing the details of one embodiment of

the local storage 106. The local storage 106 includes a microprocessor 201 (which is available from, for example, Intel Corporation or Motorola Corporation), a local (internal) bus 202, a SAN interface (I/F) 203 coupled to the SAN 107 (Figure 1), a LAN I/F 204 coupled to the LAN 108 (Figure 1), a WAN I/F 205 coupled to the WAN 110 (Figure 1), a memory I/F 206, a memory 207, a disk controller 208, and a local disk 209. The memory 207 contains the configuration table 112, statistical information 113, access log 114, cache information 115, and programs 116. The microprocessor 201 controls all resources in the local storage 106 and executes all procedures in the local storage 106 by using the programs 116.

Figure 2 shows a snapshot of the local disk 209 during the running of procedures in the local storage 106. As known to those skilled in the art, a snapshot is a copy of a whole volume or a part of a volume of a real storage. These information and programs shown in Figure 2 are stored in the local disk 209. At the boot sequence of local storage 106, these data and programs are moved from the local disk 209 to the memory 207.

The cache information 115 is configured with cache directory 210, pending data list 211, and cache data 212. The cache directory 210 is a directory (configuration information) of the cached data 212 in memory 207 and local disk 209. This cache data 212 may be managed by the least recently used (LRU) management. But some cases should be an exception of the LRU management, because there will be no possibility to access again in the near future. One such exception is sequential access where the storage address is accessed in a sequential manner.

The pending data list 211 is a list of pending data which is saved in the local storage 106. The pending data list 211 has information of each pending data such as a pointer to pending data, a pointer to a service provider site, and a synchronous period. A synchronous period is def

ined by the data feature. If the data is not important, then this data may be stored in the local storage 106 for a predetermined period. This period may be, for example, one minute, one hour, one day, one week, or one month. If the data is important, then after the data is stored in the virtual storage 111, the data should be stored immediately (synchronous) in a service provider site (e.g., site 102) by using the configuration table 112 (Figure 4). For example, synchronous data is stored in Volume 00 117. In Figure 4, Volume 00 117 is identified with the identification number 1D 00. Since the data stored in Volume 00 117 is synchronous data, the data will be immediately stored in the service provider site 102 (Remote Site 1) at Volume 10 121 (Remote 1D 10).

This synchronous access feature provides a reliable data storage system. However, synchronous access requires much time because of the access time by the user site 101 to the service provider sites 102 and/or 103. The required access time depends on the particular operation being carried out on the system 100.

Data that are accessed by hosts 104 and/or 105 are stored in the memory 207. This data is called cache data 212. Cache data 212 is a temporary data that is frequently accessed by a user of host 104 and/or host 105.

In one embodiment, the programs 116 include server programs 213, a synchronizer 217, a cache control 218, a scheduler 220, and a device driver 219. The server programs 213 permit the local storage 106 and the host 104 (and/or host 105) to communicate with each other. The server programs 213 are configured with an NFS server 214, a CIFS server 215, and an SCSI server 216. The NFS server 214 acts as a conventional NFS server. The CIFS server 215 acts as a conventional CIFS server. The SCSI server 216 acts as a SCSI target device. The synchronizer 217 synchronizes the local storage 106 and the service provider sites 102 and/or 103 by

using the pending data list 211. The cache control 218 controls cache 212 in the memory 207 and local disk 209 by using the cache directory 210. The scheduler 220 schedules all procedures such as the server programs 213, synchronizer 217, cache control 218, and device driver 219. The device driver 219 controls all devices in local storage 106 such as the SAN I/F 203, the LAN I/F 204, the WAN I/F 205, the memory I/F 206, and the disk controller 208.

Figure 3 is a block diagram showing one example of a data layout in accordance with an embodiment of the present invention. The management of the local storage 106 is based upon the use of data feature, as described below in additional details. The virtual volume storage 111 deals with the data feature, depending on whether the data is static or cached, synchronous or asynchronous. As also described further below, the virtual volume storage 111 can deal with data feature not only by storage volumes, but also by directory, file, cylinder, and/or block address. At user site 101, there are four (4) types of virtual volumes (Volume 00 117, Volume 01 118, Volume 02 119, and Volume 03 120) in the virtual storage 111. The features of these virtual volumes are described below: (1) the static synchronous volume, (2) the static asynchronous volume, (3) the cached synchronous volume, and (4) the cached asynchronous volume.

(1) The static synchronous volume (Volume 00 117):

Volume 00 117 is a static synchronous volume. Thus, Volume 00 117 occupies the same size as the real storage volume, Volume 10 121, in the first service provider site 102. The hosts 104 and/or 105 (at user site 101) read the data from Volume 00 117 and not from Volume 10 121. If the accessed data is in the cache data cache 212 (Figure 2), then the local storage 106 (Figure 2) does not access the local disk 209 (Figure 2).

and just accesses the memory 207. The hosts 104 and/or 105 write data to both Volume 00 117 and Volume 10 121. This static synchronous volume, Volume 00 117, provides a fast read access and reliable storage system.

(2) The static asynchronous Volume (Volume 01 118):

Volume 01 118 is a static asynchronous volume. Thus, Volume 001 118 occupies the same size as the real storage volume, Volume 11 122, in the first service provider site 102. The hosts 104 and/or 105 read data from just only Volume 01 118 and not from Volume 11 122. If the accessed data is in the cache data 212, then local storage 106 does not access the local disk 209 and only accesses the memory 207. The hosts 104 and/or 105 write data to Volume 01 118 and register to the pending data list 211 (Figure 2). This pending data is synchronized with a predetermined schedule, as described below.

This pending data may be compressed in the background job in the local storage 106 (Figure 2) to reduce the WAN 110 traffic. This static asynchronous volume, Volume 01 118, provides a fast read and write volume, but may not provide the same reliability as a synchronous volume.

(3) The cached synchronous volume (Volume 02 119):

Volume 02 119 is a cached synchronous volume. Thus, Volume 02 119 occupies same size as the real storage volume, Volume 12 123, in the first service provider site 102. Volume 02 119 acts as a cache. Thus, there is no data in Volume 02 119 that is accessed by the hosts 104 and/or 105. If the data is not in the local disk 209 (Figure 2), then the local storage 106 (Figure 2) reads data from the first service provider site 102 and writes the read data to the local disk 209. After the data is in the local storage 106, the local storage 106 sends the data from the local disk 209 via memory 207 to the hosts 104 and/or 105. If the accessed

d data is in the cache data 212, then the local storage 106 does not access local disk 209 and just accesses the memory 207. The hosts 104 and/or 105 write data to Volume 02 119 and Volume 12 123. This cached synchronous volume, Volume 02 119, advantageously permits the reduction of size of the local storage 106. This virtual volume, Volume 02 119, also provides a reliable storage system.

(4) The cached asynchronous volume (Volume 03 120):

Volume 03 120 is a cached asynchronous volume. Thus, Volume 03 120 occupies same size as the real storage volume, Volume 13 124, in the first service provider site 102. Volume 04 120 acts as a cache. Thus there is no data in Volume 03 120 that is accessed by the hosts 104 and/or 105. If the data is not in the local disk 209, the local storage 106 reads data from the first service provider site 102 and writes this read data into the local disk 209. After the data is already in the local storage 106, the local storage 106 sends the data from local disk 209 via memory 207 to the hosts 104 and/or 105. If the accessed data is in the data cache 207, then the local storage 106 does not access local disk 209 and just accesses the memory 207. The hosts 104 and/or 105 write the data to Volume 01 118 and register to pending data list 211 (Figure 2). This pending data is synchronized with a schedule. This cached asynchronous volume, Volume 03 120, advantageously permits the reduction of size of the local storage 106, but may not provide the same reliability as a synchronous volume.

Figure 3 also illustrates a service provided by the service provider site. Volume 11 122 in the first service provider site 102 is duplicated on volume 21 126 in the second service provider site 103. Volume 12 123 is duplicated on volume 22 127 in the second service provider site 103. Volume 13 124 is duplicated on volume 23 128 in the second service

provider site 103. This duplication (remote copying 130) of volume data provides a disaster recovery method. The volume data may be duplicated by use of, for example, conventional remote mirroring technology. Additionally, at the first service provider site 102, a backup of the data in Volume 12 123 may be made (see arrow 132) without requiring the function of the user site 101. This backup service permits the user to reduce the workload of data back up.

As now discussed below, the virtual volume storage 111 can make a trace data based upon access activities from the hosts 104 and/or 105. The virtual volume storage 111 can provide to the service provider this trace data for purposes of analyzing the user access pattern. After analyzing the user access pattern, the service provider can advantageously tune the virtual volume storage 111.

Figure 4 is a block diagram of an example of a configuration table 112 in accordance with an embodiment of the present invention. The configuration table 112 is configured with the following parameters: "ID", "Interface", "Sub Area", "Remote Site", "Remote ID", "Volume Type", "Service" ("Backup", "Analyze"), "size" ("Total", "Used", "Free"), and "Policy". ID indicates volume identification of a virtual volume in the local storage 106. For example, Volume 00 117 (Figure 3) has an ID 00. Interface indicates the virtual volume interface type such as SCSI, NFS, and CIFS. Sub area indicates the number of sub areas in the virtual volume.

A sub area may be a certain area of block address, cylinders, file, and/or directory. Volume type may be different between each sub area. In this case in the example of Figure 4, each sub area has the same volume type. For example, Volume 00 117 (Figure 3) (with ID 00) is a static synchronous volume type as mentioned previously. Remote Site indicates an identifier of the remote site (e.g., the identifier of service provider

site 102). This may be the IP address of an http address. Remote ID indicates volume ID in a service provider site. For example, Volume 10 12 1 in service provider site 102 has a remote ID 10. Volume type indicates a type of volume such as static or cached, synchronous or asynchronous. Service indicates a service provided by the service provider. Such as a back up service or an analyze service. When a service provider provides an analyze service, the service provider sets a pointer to a statistical information 113 (Figure 2) and creates statistical information (access log) 114 to obtain a data feature. If an analyze service is provided, then "Ptr" will be indicated in the Analyze section in the Configuration Table 112 as shown in the example of Figure 4. For example, Volume 02 119 (ID 02) in Configuration Table 112 has a pointer Ptr in the Analyze section (see row 400). This pointer for ID 02 (of Volume 02 119) points to ID 02 in the Statistical Information 113 in Figure 5. ID 02 has a Sub area (directory) "/usrA" (see row 500) with a pointer ("ptr") in the Log Ptr section. This pointer (ptr) points to the access log 114 which is an access log created for Volume 02 119 in the example of Figure 6.

It is further noted that in the example of Figure 5, the ID 02 (of Volume 02 119) includes a directory "/usrh" (see row 505) with a "Null" value (no pointer) in the Log Ptr section. Additionally, ID 02 includes a directory "/usrc" (see row 510) with a Null value (no pointer) in the Log Ptr section. Thus, no associated access logs 114 have been created for the directories /usrb and /usrc since these directories do not have pointers that point to an access log.

In the Configuration Table 112, Size indicates total memory size (Total) in gigabytes, used size (Used), and free size (Free). If a volume is not in the file system in local storage 106, then only the total size needs to be known. The parameter policy indicates a schedule for synchronization. For example, in row 410, the Volume 00 117 (with ID 00) has

a "null" value set in the Policy parameter in Figure 4, and this null value indicates that a synchronization schedule has not been set for Volume 00 117. Volume 02 119 (with ID 02) has a "week" value set in the Policy parameter, and this value indicates that the data in Volume 02 119 is synchronized once per week with data in a service provider site (e.g., site 102). Volume 03 120 (with ID 03) has a "day" value set in the Policy parameter, and this value indicates that the data in Volume 03 120 is synchronized daily with data in a service provider site (e.g., site 102).

The parameters in the Configuration Table 112 are set by setting the configuration data in the Configuration Table 112. Each storage subsystem has an associated Configuration Table.

Figure 5 is a block diagram showing an example of statistical information 113 in accordance with an embodiment of the present invention. When a user or service provider sets the Service-Analyze parameter to "Ptr" in the Configuration Table 112 (Figure 4), then the following will be enabled. The local storage 106 collects the user statistical access information of local storage 106 in this table of Statistical Information 113. The Statistical Information is configured with "ID", "Sub Area", "IO/s" and "MB/s" (Read Ave, Max, Write Ave, Max), "Hit Ratio", and "Log Ptr". The parameters ID and Sub area (in Figure 5) have the same meanings as ID and Sub area in Configuration Table 112 (Figure 4). IO/s indicate the number of host/user input or output activities per second such as read and write commands per second. A read or write access in a volume or in a volume directory may be recorded by use of any suitable known methods.

MB/s indicates number of bytes per second. The local storage 106 collects each read and write occurrence in a volume (or in a sub area in a

volume), and averages (Ave) this occurrences. The local storage 106 also tracks the maximum occurrence (Max) of read and write occurrence in a volume or in a sub area in a volume.

Hit Ratio indicates a hit ratio of local storage. This hit ratio is, for example, only just for the local disk 209 (Figure 2). A hit ratio is defined in equation (1):

$$(1) \text{ Hit ratio} = 100\%(\#LOCAL \text{ DISK READ}/\#ALL \text{ READS})$$

where #LOCAL DISK READ is the number of reads in the local disk 209 by a host, and #ALL READS is the number of total reads by the host.

Statistical Information 113 is sent to the service provider sites 102 and/or 103 periodically. The service provider analyzes the feature of data, and from this analysis the service provider can propose better solutions to the user of a host at the user site 101.

As stated above, Log Ptr indicates a pointer to the access log 114. If the service provider wants to know more detail of user access patterns, then the service provider can set this pointer to a created accesses log 114.

Figure 6 is a block diagram showing an example of an access log 114 in accordance with an embodiment of the present invention. When the service provider wants to analyze the details of a user's access pattern, then the service provider creates an access log 114 and sets a pointer at Log Ptr (in Statistical Information 113 in Figure 5) to this access log 114. The local storage 106 collects an access log 114 of each access by a user. This access log 114 is configured with the parameters: "Date", "Time", "Command", "File ID", "Address", and "Size". Date indicates the date of an access by a user. Time indicates the time of an access. Command indicates the command type (e.g., read or write) of an access. File ID indicates the identification of a file that was accessed by this

command. If access log is for SCSI, the File ID may be null. Address and size indicates the access address and size.

Read Process

In the read procedure, the cached device (local storage 106) acts as cache. The hosts 104 and/or 105 access the local storage 106, since the local storage 106 has a large space. At first, a host (e.g., host 104 or host 105) issues a read command of NFS, CIFS, or iSCSI protocol via LAN 108 (Figure 1), or a read command of SCSI protocol via SAN 107. The local storage 106 receives the read command via LAN I/F 204 or SAN I/F 203 (Figure 2). The device driver 219 (Figure 2) handles this read command from a host and places this read command into the scheduler 220. The scheduler 220 then analyzes the read command and places the read command into a proper server (NFS server 214, CIFS server 215, or SCSI server 216) within the server programs 213 for purposes of processing the read command. Each server program 213 checks whether the data requested by the read command is in or not in the cache data 212 (Figure 2). If all data (which is requested by the read command) is in the cache data 212, then the local storage 106 returns data to the requesting host (host 104 or 105) via internal bus 202 and a network interface (SAN I/F 203 or LAN I/F 204). On the other hand, if all of or part of the data (requested by the read command) is not in cache data 212, then the requested data should be moved to cache data 212 from the local disk 209 or from the service provider sites 102 or 103.

Figure 7 is a flowchart diagram of this read process in accordance with an embodiment of the present invention. After the local storage 106 receives a read command from a host (e.g., host 104 or 105), the local storage 106 checks 701 for the volume type of the data. The scheduler 220 checks configuration table 112 for the volume type (e.g., static type

) by check a volume ID. In particular, the scheduler 220 (Figure 2) analyzes the read command to determine the volume type of the data and to determine which server program (server 214, 215, or 216) should process (handle) the read command. The scheduler 220 looks at configuration table 112 (Figure 4) for the volume type and the appropriate server program to handle the read command. If volume type is a static volume, then the local storage 106 skips the steps 702, 703, and 704 and proceeds to step 705 which is discussed below. If the volume type is a cached volume, then the local storage 106 checks 702 if the data (requested by the read command) is stored in or not stored in the cache data 212.

The local storage 106 checks 703 the cache directory 210. In particular, the appropriate server (214, 215, or 216) in the server program 213 checks whether the data requested by the read command is in or not in the local disk 209. If all data (requested by the read command) is stored in the local disk 209, then the local storage 106 skips step 704 and proceeds to step 705 which is described below. If data (request by the read command) is not in the local disk 209, then local storage 106 executes step 704.

In step 704, the local storage 106 allocates data area in the local disk 209 and reads data (requested by the read command) from a service provider site (e.g., service provider site 102). In particular, the cache control 218 uses the configuration table 112 (Figure 4) to determine the service provider site from where data should be obtained. The cache control 218 reads the data from the service provider site and stores the data in the local disk 209.

After data is stored in local disk 209, the local storage 106 will move (read) 705 the data from the local disk 209 to the cache data 212. In particular, the cache control 218 moves the data from the local disk 209 to the cache data 212. The local storage 106 then updates 706 the s

statistical information 113 (Figure 5). In particular, the scheduler 220 in the local storage 106 updates the statistical information 113. If the log pointer (Log Ptr) (Figure 5) in the statistical information 113 is set, then the local storage 106 adds the log data of this read command to the access log 114 (Figure 6). The scheduler 220 or the device driver 219 creates the statistical information 113 and access log 114. Typically, the scheduler 220 is preferred in creating the statistical information 113 and access log 114.

The data (requested by the read command) is then moved 707 from the cache data 212 to the host (e.g., host 104 or 105) that sent the read command. In particular, the cache control 218 moves the data from the cache data 212 to the appropriate server program (214, 215, or 216), and the appropriate server program (214, 215, or 216) returns the data to the host that sent the read command.

Write process

In the write procedure, the cached device (local storage 106) acts as a cache. The host 104 and/or host 105 access the local storage 106, since the local storage 106 has a large space. At first, a host (e.g., host 104 or host 105) issues a write command of NFS, CIFS, or iSCSI protocol via LAN 108, or a write command of SCSI protocol via SAN 107. The local storage 106 receives the write command via LAN I/F 204 or SAN I/F 203. The device driver 219 (Figure 2) handles this read command from a host and places this read command into the scheduler 220. The scheduler 220 then analyzes the write command and places the write command into a proper server (NFS server 214, CIFS 215, or SCSI server 216) within the server programs 213. Each server program 213 checks whether the data of the write command is in or not in the cache data 212 (Figure 2). If all data (of the write command) is in the cache data 212, then the local s

storage 106 receives the data from the host via internal bus 202 and a network interface (SAN I/F 203 or LAN I/F 204). The local storage 106 will overwrite the data on an allocated space in the local disk 209 and cache 212. On the other hand, if all of or part of the data (of the write command) is not in the cache data 212, then the local storage 106 will allocate an area (in local disk 209 and cache data 212) for the rest of data not in the cache data 212. After all of the data area is allocated, then the local storage 106 stores data in that allocated area. When all data is stored in the cache data 212, then the local storage 106 stores data in local disk 209 and sends this data to a provider site (e.g., service provider site 102).

Figure 8 is a flowchart diagram of a write process in accordance with an embodiment of the present invention. After the local storage 106 receives a write command from a host, the local storage 106 checks 801 the volume type of the data of the write command. In particular, the scheduler 220 (Figure 2) analyzes the write command to determine the volume type of the data and to determine which server program (server 214, 215, or 216) should process (handle) the read command. The scheduler 220 looks at configuration table 112 (Figure 4) for the volume type and the appropriate server program to handle the write command. If the volume type is a static volume, then the local storage 106 skips steps 802, 803, and 804, and proceeds to step 805 which is described below. If the volume type is a cached volume, then the local storage 106 checks if the data (of the write command) is stored or not stored in cache data 212 (Figure 2).

The local storage 106 checks 802 the cache directory 210 (Figure 2).

In particular, the appropriate server (214, 215, or 216) in the server program 213 checks whether the data of the write command is allocated in the local disk 209. In step 803, if all the data (requested by the wr

ite command) is allocated in the local disk 209, then the local storage 106 skips step 804 and proceeds to step 805. In step 803, if the data is not in the local disk 209, then local storage 106 executes step 804.

If the data (of the write command) is not allocated in the local disk 209, then the local storage 106 allocates 804 data area in both the local disk 209 and cache data 212 for the write data. In particular, the cache control 218 allocates the data area on both the local disk 209 and cache data 212.

The local storage 106 then writes 805 the write data to the local disk 209 via the cache data 212. In particular, the appropriate server program (NFS server 214, CIFS server 215, or SCSI server 216) writes the data to the local disk 209.

The local storage 106 checks 806 the volume type in which the data is written by use of the configuration table 112 (Figure 4). In particular, the server program (NFS server 214, CIFS server 215, or SCSI server 216) uses the configuration table 112 to check for the volume type. If the volume type is a synchronous area, then the local storage 106 immediately synchronizes (writes) 808 the data to a service provider site (e.g., site 102) and proceeds to step 809 which is discussed below. The synchronizer 217 (Figure 2) performs the synchronizing of the data (in the user site) to the service provider site. The cache control 218 in the local storage 106 knows where the data should be stored in the service provider site by use of the configuration table 112.

In step 806, if the volume type is an asynchronous area, then the local storage 106 updates 807 the pending data list 211 (Figure 2). In particular, the server program (server 214, 215, or 216) updates the pending data list.

The local storage 106 then updates 809 the statistical information 113 (Figure 5). In particular, the scheduler 220 updates the statistical

information 113. If the log pointer (Log Ptr) (Figure 5) in statistical information 113 is set, then the local storage 106 adds the log data of this write command to the access log 114. In particular, the scheduler 220 adds the log data of the write command to the access log 114. The write process then ends.

Synchronization Process

In the case of an asynchronous volume write (a data write to an asynchronous area), the data should be sent to a service provider site (e.g., site 102) by a schedule that is defined by the user or service provider. This schedule is, for example, defined by the value in the "Policy" entry in the configuration table 112 in Figure 4. The scheduler 220 (Figure 2) executes the synchronizer 217 (Figure 2) periodically to enable data synchronization with a service provider site. This period may be set to, for example, approximately 1.0 milli-second or 10.0 milli-seconds, depending on the system 100 setting.

Figure 9 is a flowchart diagram of synchronization process in accordance with an embodiment of the present invention. The synchronization process may be performed by the synchronizer 217 (Figure 2).

The synchronizer 217 selects 901 the head data in the pending data list 211 (Figure 2). At the first, the synchronizer 217 checks 902 the pending data list 211. If there is no data in the pending data list 211, then the synchronizer 217 ends the synchronization process. If there is one or more data (in the pending data list 211) which should be synchronized, then the synchronizer 217 executes steps 902-905 as described below.

The synchronizer 217 checks 903 the information of pending data list 211. If this pending data in the pending data list 211 should be synchronized, then the synchronizer 217 sends the pending data to a service p

provider site (e.g., site 102), so that the pending data is synchronized 904 with the service provider site. The cache control 218 in the local storage 106 knows where data should be stored in the service provider site (e.g., site 102) based upon the values in the Remote ID and Remote Site in the configuration table 112 (Figure 4).

The synchronizer 217 then selects 905 the next data in the pending data list 211. The synchronizer 217 repeats steps 902 through 905 until there is no more data to select in the pending data list 211. If there is no more data to select in the pending data list 211, then the method of Figure 9 ends.

It is also within the scope of the present invention to implement a program or code that can be stored in an electronically-readable medium to permit a computer to perform any of the methods described above.

Thus, while the present invention has been described herein with reference to particular embodiments thereof, a latitude of modification, various changes and substitutions are intended in the foregoing disclosure, and it will be appreciated that in some instances some features of the invention will be employed without a corresponding use of other features without departing from the scope of the invention as set forth.

4 Brief Description of Drawings

Figure 1 is block diagram of a system in accordance with an embodiment of the present invention;

Figure 2 is a block diagram showing additional details of one embodiment of the local storage in Figure 1;

Figure 3 is a block diagram showing one example of a data layout in accordance with an embodiment of the present invention;

Figure 4 is a block diagram of an example of a configuration table in accordance with an embodiment of the present invention;

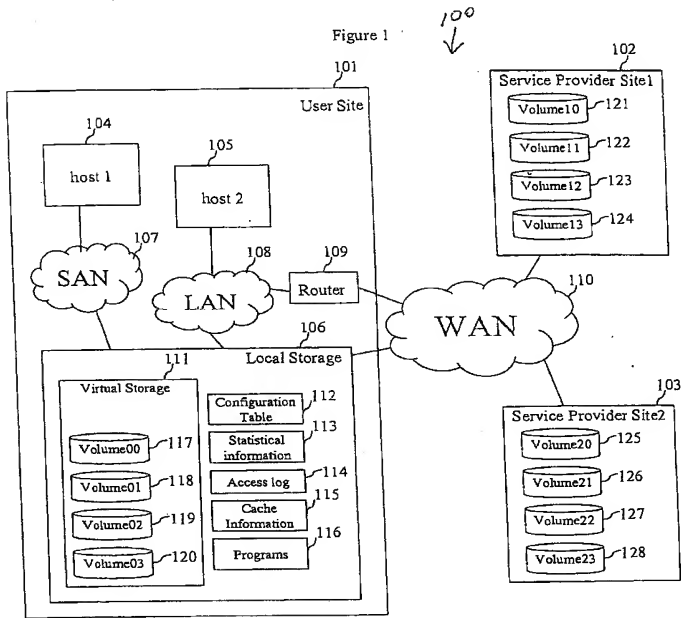
Figure 5 is a block diagram of an example of statistical information in accordance with an embodiment of the present invention;

Figure 6 is a block diagram of an example of an access log in accordance with an embodiment of the present invention;

Figure 7 is a flowchart diagram of a read process in accordance with an embodiment of the present invention;

Figure 8 is a flowchart diagram of a write process in accordance with an embodiment of the present invention; and

Figure 9 is a flowchart diagram of a synchronization process in accordance with an embodiment of the present invention.



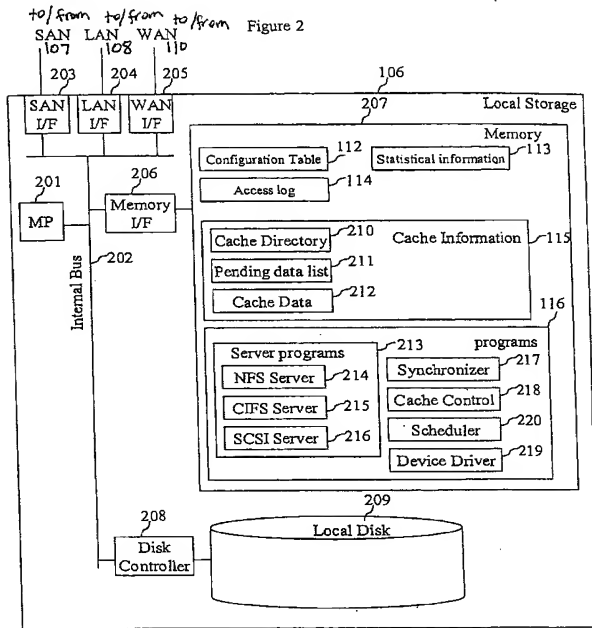


Figure 5

Statistical information

		IO/s				MB/s				Hit Ratio	Log Ptr
ID	Sub Area	Read		Write		Read		Write			
		Ave	Max	Ave	Max	Ave	Max	Ave	Max		
00											
01		15	100	3	25	0.12	2.0	0.02	0.5	80%	Null
02	/usrc										ptr
	/usrb										Null
	/usrc										Null
⋮											

113

500

505

505

Figure 6

Access log 114

Date	Time	Command	File ID	Address	Size
2000.8.1	0:0:0	Read	0	0x0000	0x0001
2000.8.1	0:0:0	Write	0	0x0000	0x0001
2000.8.1	0:0:1	Read	10	0x0100	0x0100
2000.8.1	0:0:1	Read	10	0x0200	0x0100
2000.8.1	0:0:2	Read	10	0x0200	0x0100
2000.8.1	0:0:2	Read	10	0x0300	0x0100
2000.8.1	0:0:2	Read	10	0x0400	0x0100
⋮					

Figure 7

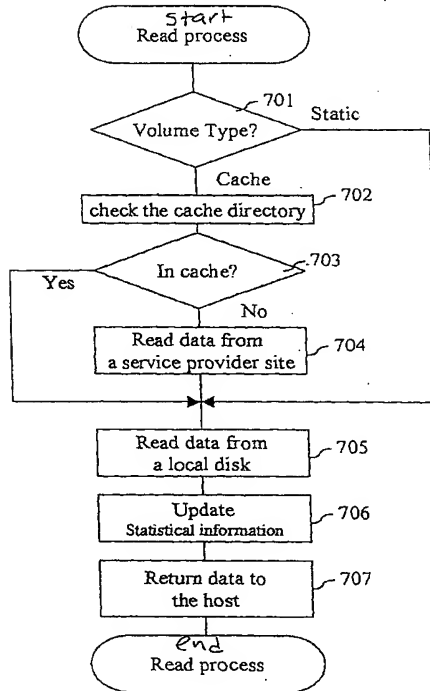


Figure 8

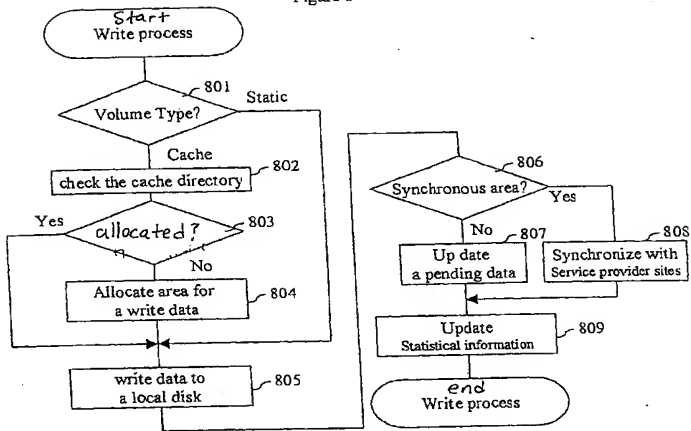
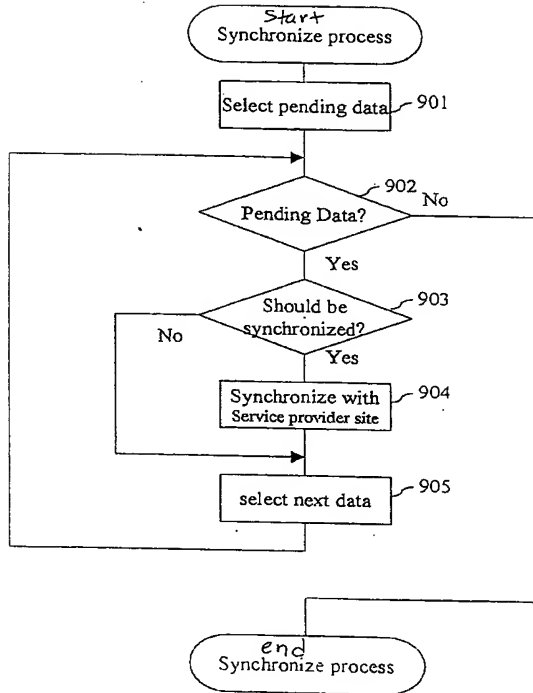


Figure 9



1. Abstract

A system for providing a data storage service, comprises: a service provider site configured to provide a data storage service; and a user site coupled by a wide area network (WAN) to the service provider site, the user site comprising a local storage having a virtual storage, the virtual storage having a synchronous volume and an asynchronous volume, the local storage configured to immediately transmit to the service provider site data that is written in the synchronous volume, to transmit at a predetermined schedule to the service provider site data that is written in the asynchronous volume, and to read data from the service provider site if the data is not stored in the local storage.

2. Representative Drawing

Fig. 1